

CCT College Dublin

Assessment Cover Page

Module Title:	Strategy Thinking
Assessment Title:	Project Report
Lecturer Name:	James Garza
Student Full Name:	Giulio Calef, Kevin Byrne and Victor Ferreira Silva
Student Number:	sba22314, sba22264, 2021324
Assessment Due Date:	7th May 2023
Date of Submission:	7th May 2023

Declaration

By submitting this assessment, I confirm that I have read the CCT policy on Academic Misconduct and understand the implications of submitting work that is not my own or does not appropriately reference material taken from a third party or other source. I declare it to be my own work and that all material from third parties has been appropriately referenced. I further confirm that this work has not previously been submitted for assessment by myself or someone else in CCT College Dublin or any other higher education institution.

Detecting and Predicting Severe Slugging in Petrobras 3W Data Set

Strategic Thinking Capstone Project

**Giulio Calef
Kevin Byrne
Victor Ferreira Silva**

Strategic Thinking Capstone Project

Higher Diploma in Science in Artificial Intelligence Applications
CCT College Dublin
Ireland
May 2023

Contents

1	Introduction	2
1.1	Business Understanding	2
1.2	Hypothesis	2
1.3	General Goal	2
1.4	Success Criteria	2
1.5	Methodologies and Technologies	3
1.6	Accomplishment	4
2	Data Understanding	4
2.1	Data Characterisation	4
2.2	Exploratory Data Analysis	6
3	Data Preparation	6
3.1	Data Cleaning	6
3.2	Feature Engineering	8
3.3	Train/Test Splitting	9
3.4	Handling Imbalanced Data	10
3.5	Data Scaling	10
3.6	Dimensionality Reduction	10
4	Modeling	10
4.1	Baseline: DummyClassifier	12
4.2	LinearSVC	12
4.2.1	Hyperparameter optimisation	12
4.2.2	Model training	13
4.3	k-Nearest Neighbors Classifier	13
4.3.1	Hyperparameter optimisation	13
4.3.2	Model training	14
4.4	Artificial Neural Network	14
4.4.1	Model definition	14
4.4.2	Hyperparameter optimisation	15
4.4.3	Model training	16
4.5	Decision Tree Classifier	17
4.5.1	Hyperparameter optimisation	17
4.5.2	Model training	18
4.6	Random Forest Classifier	18
4.6.1	Hyperparameter optimisation	18
4.6.2	Model training	19
5	Evaluation	19
5.1	Classification Report	19
5.1.1	LinearSVC	19
5.1.2	k-Neighbors Classifier	19
5.1.3	Neural Networks	20
5.1.4	Decision Tree	20
5.1.5	Random Forest	20
5.2	Confusion Matrices	20
5.3	10-Fold Cross Validation	22
6	Conclusion	23

Detecting and Predicting Severe Slugging in Petrobras 3W Data Set

Giulio Calef, Kevin Byrne, Victor Ferreira Silva

May 4, 2023

1 Introduction

1.1 Business Understanding

Operational safety, productivity, quality are general key objectives in any industry, and in the oil and gas industry, environmental concerns not only are crucial importance, but poses as an immense daily challenge. Any events which can cause a production loss in an oil field are certainly very costly, in pure economic terms, but can also have dire costs in terms of human lives and of damage to our environment that can be hard to address or almost permanent in some cases. Therefore, any technique that can help early detection and/or prevention of technical accidents in the oil industry is therefore very welcome and “worth as gold”. Offshore oil wells provide some of the most challenging operating conditions in the industry, with additional complexity due to the peculiarities of operating at sea and the limited amount of instrumentation that can be deployed to monitor and control the well operational status.

A simplified description of a typical offshore well can be seen below in Figure 1.

Basically, a key structure in the well is the so-called *Christmas Tree*, a structure lying on the seabed, at the well head, with pressure and temperature sensors (PDG – Permanent Downhole Gauge and TPT – Temperature and Pressure Transducer), and safety valves (DHSV – Downhole Safety Valve). In turn, a PCK (Production Choke valve) is installed on the drilling vessel/rig at the top.

Given this, Petrobras, the Brazilian oil company, has developed a data set (3W) that contains data for the most common monitored variables in offshore oil wells. This project aims to understand the relations between these variables and a specific oil well operational anomaly named *Severe Slugging*.

Severe Slugging is an critical flow assurance issue, commonly observed in offshore pipeline-riser systems. Some of the consequences of this issue include flooding of downstream production facilities and an overall decrease in productivity.

1.2 Hypothesis

The data present in 3W Data Set allows classifier models to predict Severe Slugging with high accuracy.

1.3 General Goal

The business objective of our project is to apply machine learning techniques modeled in this project to the 3W data set to accurately predict *Severe Slugging* in an offshore well production line. Achieving a significant accuracy in these predictions should be possible by identifying the correlations between the variables presented in the records monitored in an offshore oil well operation.

1.4 Success Criteria

Success will be defined as the ability to accurately predict one or more of the conditions leading to an undesirable event. Given the potentially catastrophic impacts (ethical, environmental, economic) of accidents in the oil industry, the ability to predict a potential risk so that a quick reaction/fix can avoid unrecoverable conditions has a clear business value.

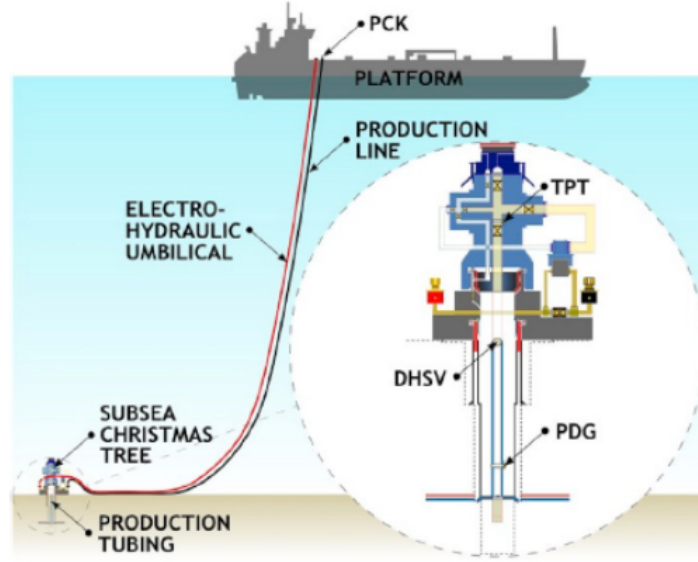


Figure 1: Schematic of a typical offshore well - source: Petrobras

1.5 Methodologies and Technologies

The methodology that guided the iterations this project was CRISP-DM. Given this, a number of libraries were used throughout this process to support most of its distinct stages, that is, Data Understanding, Data Preparation, Modeling, and Evaluation.

Once the methodology was defined, Petrobras' *3W Tool Kit* was studied and used to extract the data from the instances interesting to the project, that is, from the real instances that effectively presented the Severe Slugging event.

Then, *Pandas* and *NumPy* libraries were selected for data manipulation and analysis. Then, the base machine learning library for the majority of models in this project was *Scikit-learn*, as it offers various preprocessing, classification models and clustering algorithms. Another library imported in this project was *Keras*, which is an open-source solution that provides an interface for artificial neural networks.

Besides that, *StandardScaler* from the *Scikit-learn* library was selected for scaling and normalisation of the data. Given the high data imbalance presented by the data set, a *RandomUnderSampler* was also from *Imbalanced-learn* library was imported. Regarding feature reduction, the decomposition algorithm *PCA*, also from *Scikit-learn*, was adopted here in this study.

Also, *Seaborn*, *Matplotlib* and *Plotly* libraries were widely used to analyse data throughout the project and provide data visualisation using at the Evaluation stage. Another tool from *Scikit-learn* selected as a baseline model was *DummyClassifier*.

Additionally, the following tools from module *model_selection* in *Scikit-learn* library were chose:

- *GridSearchCV*, for hyperparameter tuning
- *cross_val_score* for cross-validation,
- *train_test_split* for train/test splitting,
- *KFold* for cross-validation during Evaluation

Ultimately, *Scikit-learn* library also provided the following classifiers for this project:

- *LinearSVC*, from *svm* module,
- *KNeighborsClassifier*, from *neighbors* module,
- *DecisionTreeClassifier*, from *tree* module,

- RandomForestClassifier, from *ensemble* module

Lastly, the modules *make_pipeline* and *Pipeline* from *Scikit-learn* were used to chain all steps of the workflow together, and the library *Pickle* was selected to persist the resulting models in file.

1.6 Accomplishment

After extracting using 3W Tool Kit and some CRISP-DM iterations, this project presented 2 models with a high accuracy.

2 Data Understanding

Preprocessing a dataset through data characterisation involves summarising the features and characteristics present in the data using statistical measures and visualisations techniques such as bar charts and scatter plots. After this stage, it should be possible to identify biases, patterns, trends, and any missing or irrelevant data in the data set that may need to be addressed.

This dataset is composed by instances of eight types of undesirable events characterized by eight process variables from three different sources: real instances, simulated instances and hand-drawn instances. All real instances were taken from the plant information system that is used to monitor the industrial processes at an operational unit in Brazilian state of Espírito Santo. The simulated instances were all generated using <https://www.software.slb.com/products/olga>, a dynamic multiphase flow simulator that is widely used by oil companies worldwide (Andreolli, 2016). Finally, the hand-drawn instances were generated by a specific tool developed by Petrobras researchers for this dataset to incorporate undesirable events classified as rare.

Ultimately, only the data from the real instances were select for this project, as simulated data and hand-drawn instances did not present any record for two relevant features, namely Gas Lift Flow Rate and Pressure Variable Upstream Of the Gas Lift Choke.

2.1 Data Characterisation

The data consists of over 50 million observations, with 13 columns of data for each observation. The first column, label, indicates the event type for each observation. The second column, well, contains the name of the well the observation was taken from. Hand-drawn and simulated instances have fixed names for in this column, while real instances have names masked with incremental id. The third column, id, is an identifier for the observation and it is incremental for hand-drawn and simulated instances, while each real instance has an id generated from its first timestamp. The columns representing the process variables are:

- P-PDG: pressure variable at the Permanent Downhole Gauge (PDG) - installed on Christmas Tree;
- P-TPT: pressure variable at the Temperature and Pressure Transducer (TPT) - installed on Christmas Tree;
- T-TPT: temperature variable at the Temperature and Pressure Transducer (TPT);
- P-MON-CKP: pressure variable upstream of the production choke (CKP) - located on platform;
- T-JUS-CKP: temperature variable downstream of the production choke (CKP);
- P-JUS-CKGL: pressure variable upstream of the gas lift choke (CKGL);
- T-JUS-CKGL: temperature variable upstream of the gas lift choke (CKGL);
- QGL: gas lift flow rate;

The pressure features are measured in Pascal (Pa), the volumetric flow rate features are measured in standard cubic meters per second (SCM/s), and the temperature features are measured in degrees Celsius (°C).

Other information are also loaded into each pandas Dataframe:

Column	pandas.Dtype
timestamp	datetime64[ns]
label	int64
well	object
id	int64
P-PDG	float64
P-TPT	float64
T-TPT	float64
P-MON-CKP	float64
T-JUS-CKP	float64
P-JUS-CKGL	float64
T-JUS-CKGL	float64
QGL	float64
class	float64
source	object

Table 1: Summary of the data set compiled from real instances

- label: instance label (event type) - target variable;
- well: well name. Hand-drawn and simulated instances have fixed names (respectively, drawn and simulated. Real instances have names masked with incremental id;
- id: instance identifier. Hand-drawn and simulated instances have incremental id. Each real instance has an id generated from its first timestamp;
- class: Although it can be used to identify periods of normal operation, fault transients, and faulty steady states, which can help with diagnosis and maintenance, it is a category which results from label, which is our target here

The labels are:

- 0 - Normal Operation = Normal
- 1 - Abrupt Increase of BSW = AbrIncrBSW
- 2 - Spurious Closure of DHSV = SpurClosDHSW
- 3 - Severe Slugging = SevSlug
- 4 - Flow Instability = FlowInst
- 5 - Rapid Productivity Loss = RProdLoss
- 6 - Quick Restriction in PCK = QuiRestrPCK
- 7 - Scaling in PCK = ScalingPCK
- 8 - Hydrate in Production Line = HydrProdLine

In order to maintain the realistic aspects of the data, the data set was built without pre-processing, including the presence of NaN values, frozen variables due to sensor or communication issues, instances with varying sizes, and outliers (Vargas et al. 2019).

A concise summary of this data set generated by *pandas.DataFrame.info* method can be seen on Table 1.

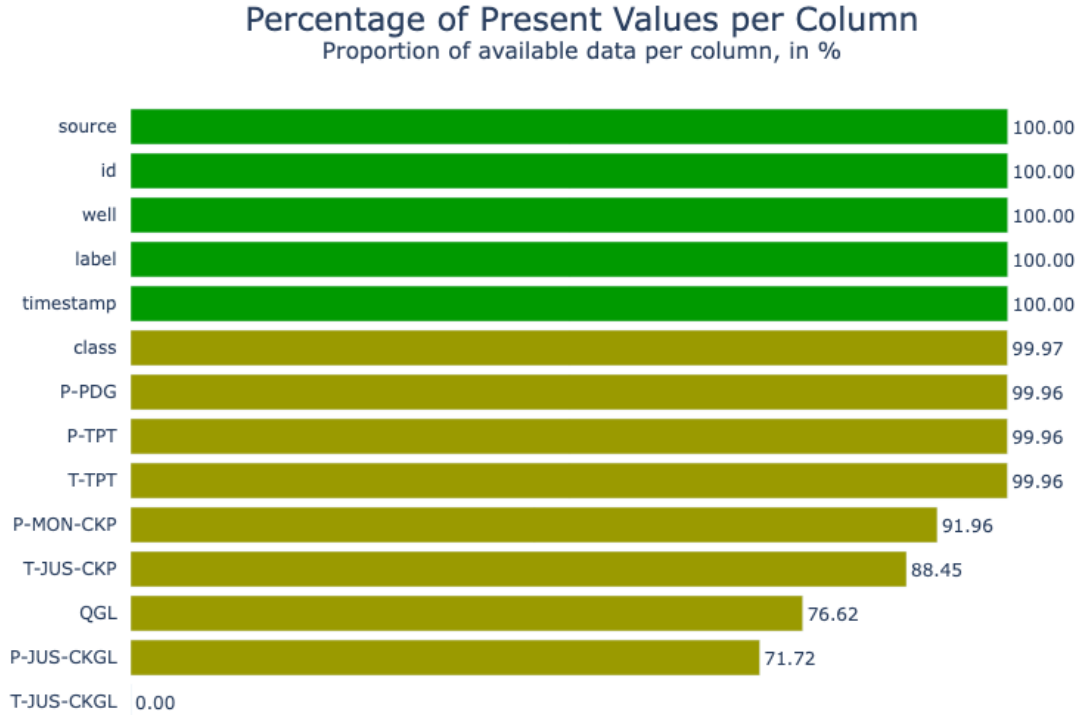


Figure 2: Proportion of available data per column, in %.

2.2 Exploratory Data Analysis

A bar chart was generated displaying the percentage of present values in each column of the data frame - see Figure 2. It contained missing values in several columns, thus some columns and row were deleted in order to obtain accurate and reliable results.

Three boxplots were plotted to show how the data was distributed before any data cleaning - see Figure 3. They were divided according the feature measurement unit: the pressure features were measured in Pascal (Pa), the temperature features are measured in degrees Celsius ($^{\circ}\text{C}$) and one feature about volumetric flow rate which was measured in standard cubic meters per second (SCM/s).

3 Data Preparation

Data preparation included Data Cleaning, Feature Engineering, Train/Test Splitting and Handling Imbalanced Data, Data Scaling, and an analysis of the chosen approach regarding dimensionality reduction for some models.

3.1 Data Cleaning

The missing data from the following columns were removed: class, P-PDG, P-TPT, T-JUS-CKP, P-MON-CKP, T-TPT, P-MON-CKP, QGL and P-JUS-CKGL. After this, the columns class, T-JUS-CKGL (an empty column), id, source were dropped. Column class is a column which brings more details about label. Consider that columns timestamp, label were kept at this stage. Finally all duplicates were removed.

```

1 # dropping rows with missing or null class column
2 df_clean = df.dropna(subset=[

```

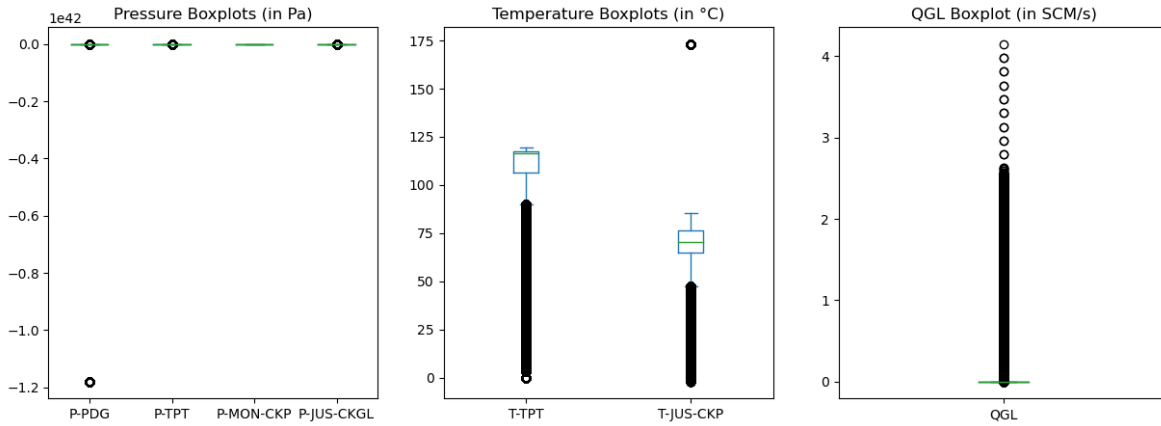



Figure 3: Box plots showing the distribution of pressure, temperature, and QGL (SCM/s) data for a set of oil wells.

```

3     'class', 'P-PDG', 'P-TPT', 'T-JUS-CKP', 'P-MON-CKP', 'T-TPT',
4     'P-MON-CKP', 'QGL', 'P-JUS-CKGL'
5 ])
6
7 # removing redundant columns
8 df_clean = df_clean.drop(['class', 'T-JUS-CKGL', 'id', 'source'], axis=1)
9
10 # checking duplicated rows after removing ids
11 df_clean = df_clean.drop_duplicates()
12
13 df_clean.info()

```

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 10003580 entries, 0 to 13952910
Data columns (total 10 columns):
#   Column      Dtype
---  -
0   timestamp   datetime64[ns]
1   label       int64
2   well        object
3   P-PDG       float64
4   P-TPT       float64
5   T-TPT       float64
6   P-MON-CKP   float64
7   T-JUS-CKP   float64
8   P-JUS-CKGL  float64
9   QGL         float64
dtypes: datetime64[ns](1), float64(7), int64(1), object(1)
memory usage: 839.5+ MB

```

Also, as it can be seen on Figure 3, features P-PDG and P-TPT had the presence of extreme outliers. These outliers were also removed with the following code:

```

1 # removing extreme outliers from P-PDG
2 Q1 = df_clean['P-PDG'].quantile(0.25)
3 Q3 = df_clean['P-PDG'].quantile(0.75)
4 IQR = Q3 - Q1
5 lower_bound = Q1 - (3 * IQR)
6 df_no_outliers = df_clean[(df_clean['P-PDG'] >= lower_bound)]
7
8 # removing extreme outliers from P-TPT
9 Q1 = df_no_outliers['P-TPT'].quantile(0.25)
10 Q3 = df_no_outliers['P-TPT'].quantile(0.75)
11 IQR = Q3 - Q1

```

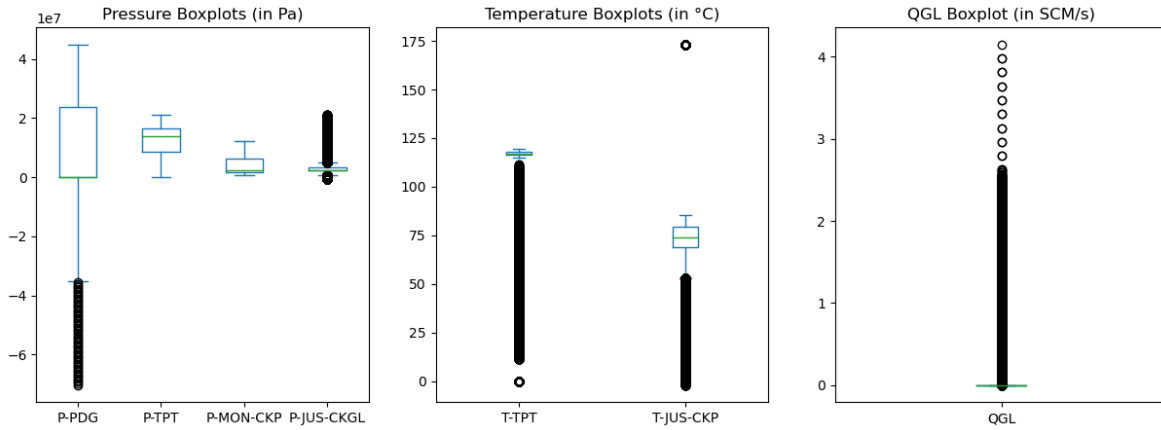


Figure 4: Box plots showing the distribution of pressure, temperature, and QGL (SCM/s) data without extreme outliers.

```

12 upper_bound = Q3 + (3 * IQR)
13 df_no_outliers = df_no_outliers[(df_no_outliers['P-TPT'] <= upper_bound)]
14
15 df_no_outliers.shape

```

(9780901, 10)

These rows with presence of extreme outliers represented 2.26% of the resulting rows so far. As a result the distribution of values in P-PDG and P-TPT were modified, as Figure 4 shows.

3.2 Feature Engineering

Given the label feature contains 8 possible numeric labels for each undesirable event and 1 label value 0 for normal observations, 8 new boolean columns were created for each one undesirable event, including for Severe Slugging, which is this project's target.

```

1 dt_feat = df_no_outliers
2
3 # Changing 'label' column to object dtype
4 dt_feat['label'] = dt_feat['label'].astype('object')
5
6 # Creating uint8 columns for each label
7 label_dummies = pd.get_dummies(dt_feat['label'], prefix='label')
8 dt_feat = pd.concat([dt_feat, label_dummies], axis=1)
9
10 # Renaming uint8 columns
11 column_names = {
12     'label_0': 'Normal',
13     'label_1': 'AbrIncrBSW',
14     'label_2': 'SpurClosDHSW',
15     'label_3': 'SevSlug', # target
16     'label_4': 'FlowInst',
17     'label_5': 'RProdLoss',
18     'label_6': 'QuiRestrPCK',
19     'label_7': 'ScalingPCK',
20     'label_8': 'HydrProdLine'
21 }
22 dt_feat = dt_feat.rename(columns=column_names)
23
24 # Dropping the original 'label' column and Normal column,
25 # since all other events must be 0
26 dt_feat = dt_feat.drop(['label', 'Normal'], axis=1)
27 dt_feat.info()

```

<class 'pandas.core.frame.DataFrame'>

```

Int64Index: 9780901 entries, 0 to 13952910
Data columns (total 16 columns):
#   Column      Dtype
---  -
0   timestamp   datetime64[ns]
1   well         object
2   P-PDG        float64
3   P-TPT        float64
4   T-TPT        float64
5   P-MON-CKP    float64
6   T-JUS-CKP    float64
7   P-JUS-CKGL   float64
8   QGL          float64
9   AbrIncrBSW   uint8
10  SpurClosDHSW uint8
11  SevSlug      uint8
12  FlowInst     uint8
13  RProdLoss    uint8
14  QuiRestrPCK  uint8
15  ScalingPCK   uint8
dtypes: datetime64[ns](1), float64(7), object(1), uint8(7)
memory usage: 811.5+ MB

```

Then all undesirable events columns were deleted but the column which denotes the observations presents Severe Slugging. The column *HydrProdLine* concerned to Hydrate in Production line, however this event was not found in the data set resulting from real instances.

```

1 dt_feat_target = dt_feat.drop([
2 #     , 'SevSlug', 'HydrProdLine',
3     'AbrIncrBSW', 'SpurClosDHSW', 'FlowInst', 'RProdLoss', 'QuiRestrPCK', 'ScalingPCK'
4 ], axis=1)
5
6 dt_feat_target.info()

```

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 9780901 entries, 0 to 13952910
Data columns (total 10 columns):
#   Column      Dtype
---  -
0   timestamp   datetime64[ns]
1   well         object
2   P-PDG        float64
3   P-TPT        float64
4   T-TPT        float64
5   P-MON-CKP    float64
6   T-JUS-CKP    float64
7   P-JUS-CKGL   float64
8   QGL          float64
9   SevSlug      uint8
dtypes: datetime64[ns](1), float64(7), object(1), uint8(1)
memory usage: 755.6+ MB

```

3.3 Train/Test Splitting

The following code defined how the data set was split in Train and Test data sets. Additionally, the columns *timestamp* and *well* were removed and at the end the percentual distribution of the records according the presence or absence of Severe Slugging was computed.

```

1 # defining features (X) and label (y)

```

```

2 target = 'SevSlug'
3
4 X = dt_feat_target.drop([target, 'timestamp', 'well'], axis=1)
5 y = dt_feat_target[target]
6
7 # splitting data into train and test sets
8 X_train_u, X_test, y_train_u, y_test = train_test_split(X, y, test_size=0.3,
9 random_state=42)
10
11 class_names = {0: 'Non Sev Slug', 1: 'SEV SLUGGING'}
12 print(y_train_u.value_counts(normalize=True).rename(index=class_names))

```

```

Non Sev Slug    0.94194
SEV SLUGGING    0.05806
Name: SevSlug, dtype: float64

```

After the splitting process, the training data set had 6,846,630 rows and the test data set had 2,934,271 rows.

3.4 Handling Imbalanced Data

A *RandomUnderSampler* was chosen to balance training data. As a result 50% of observations presented Severe Slugging while the other 50% were normal or presented other undesirable event. Test data set was not balanced since this project aim to best represent your deployment scenarios in real life.

```

1 # balancing data
2 balancing = RandomUnderSampler(random_state=42)
3
4 X_train, y_train = balancing.fit_resample(X_train_u, y_train_u)
5
6 class_names = {0: 'Non Sev Slug', 1: 'SEV SLUGGING'}
7 print(y_train.value_counts(normalize=True).rename(index=class_names))
8 print([X_train.shape, y_train.shape])

```

```

Non Sev Slug    0.5
SEV SLUGGING    0.5
Name: SevSlug, dtype: float64
[(795026, 7), (795026,)]

```

Handling data imbalance is also important because it affects correlations - see as Figure 5 shows.

3.5 Data Scaling

Although there are features presenting non-normal distributions, *StandardScaler* was chosen as data scaler. It was chose because there are some features with strong correlation with Severe Slugging and lognormal distributions such as *QGL* and *P-JUS-CKGL* and as it is a method sensitive to the presence of outliers. The results of this transformation can be seen on Figure 6.

3.6 Dimensionality Reduction

The unsupervised learning technique Principal Component Analysis (PCA) was chosen not only to prepare the data for some of the models studied here, but also to evidence any possible linear separability in this model. In Figure 7 the results of this dimensionality reduction can be seen in two ways, with 2 and 3 components, although this process was unnecessary for the most successful models, that is, the non-linear classifiers.

4 Modeling

For this project, five models were chosen: LinearSVC, k-Nearest Neighbors Classifier, Artificial Neural Network, Decision Tree Classifier and Random Forest Classifier.

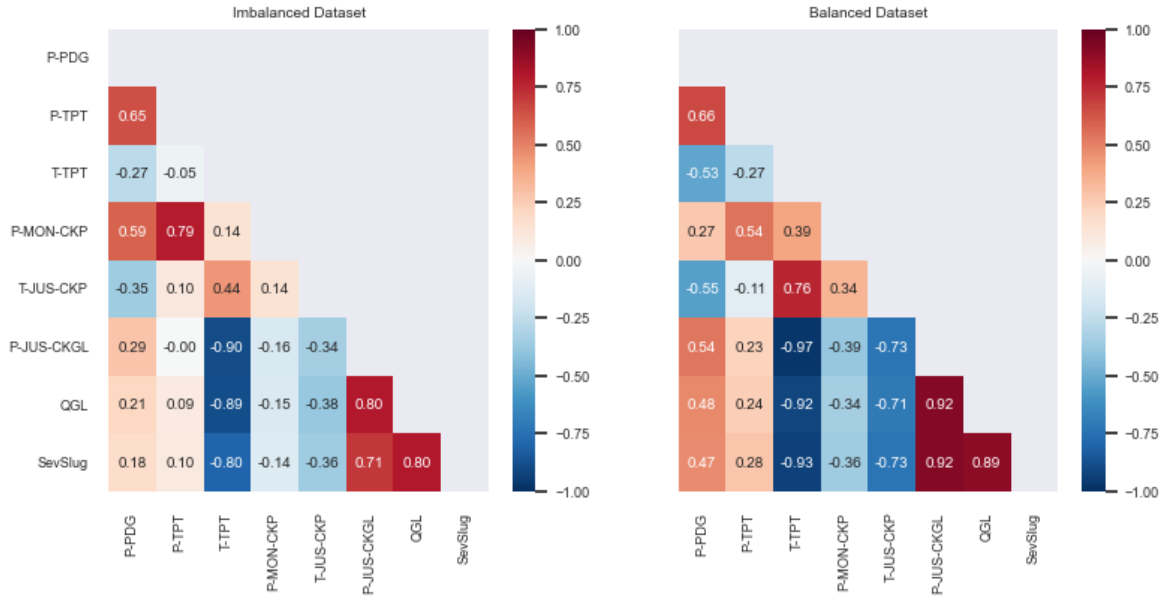


Figure 5: Correlations between variables before and after data balancing

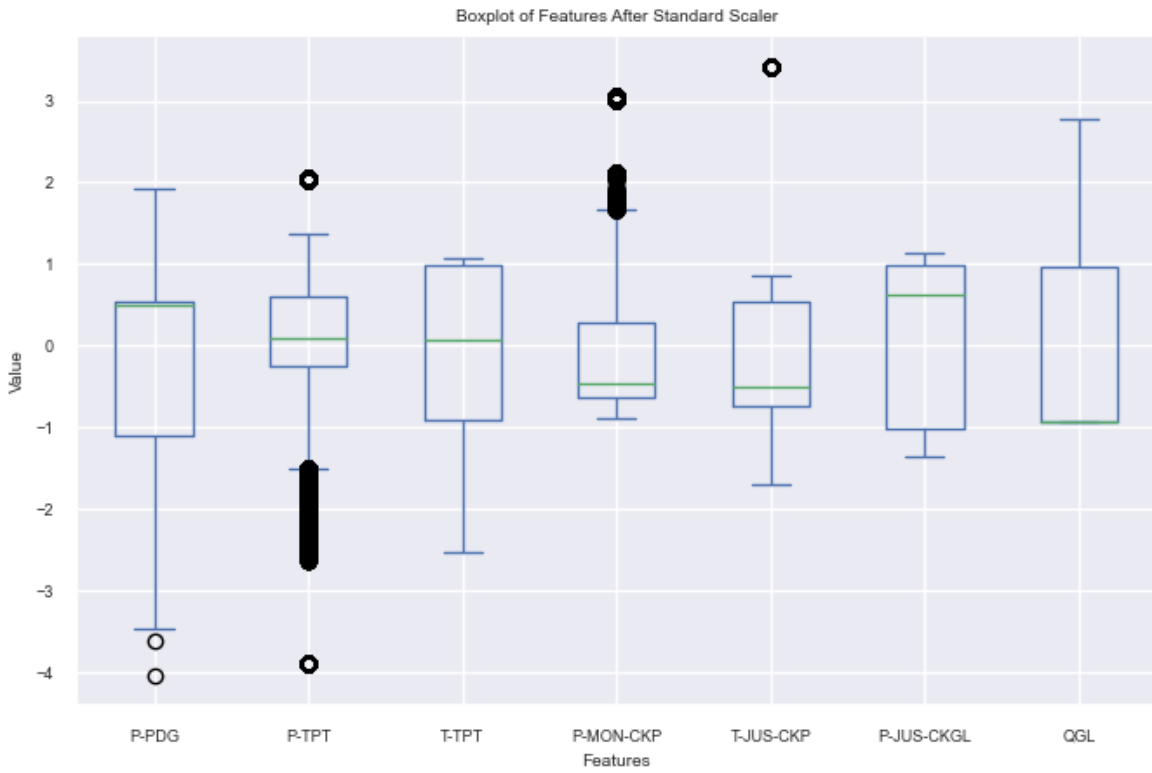


Figure 6: Box plot showing the distribution of the features in the training set after applying the StandardScaler transformation

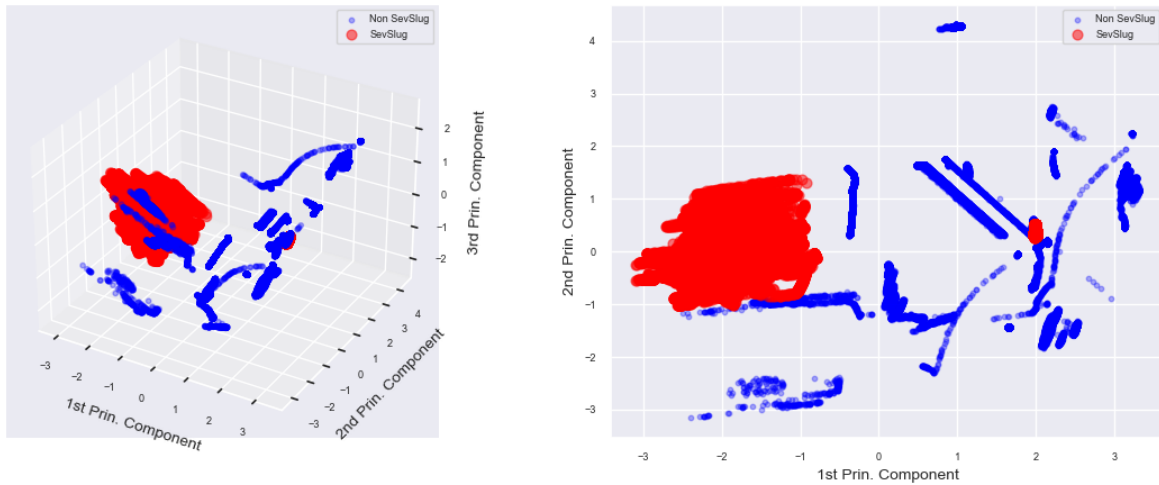


Figure 7: Visualisation of PCA applied to the data set showing a scatter plot for two (2D) and three (3D) principal components.

4.1 Baseline: DummyClassifier

Before proceeding with other models, a DummyClassifier was adopted to find a baseline for validation accuracy using the test data set. Given the distribution of the target variable in test data set, the baseline for any model was set in 94.17%.

```
1 dummy_pipeline = make_pipeline(StandardScaler(), DummyClassifier())
2 dummy_pipeline.fit(X_train, y_train)
3
4 # confirming score for Dummy classifier results from a balanced dataset
5 score = dummy_pipeline.score(X_train, y_train)
6
7 # predicting
8 y_predicted = dummy_pipeline.predict(X_test)
9 baseline = metrics.accuracy_score(y_test, y_predicted)
10
11 print("Score: ", score)
12 print("Accuracy: ",baseline)
```

Score: 0.5

Accuracy: 0.9417780429960286

4.2 LinearSVC

Although the data is not linearly separable and it's not possible to find a condition of 100% correctly classified by a hyperplane, a linear support vector classifier (LinearSVC) was implemented as part of this benchmark.

4.2.1 Hyperparameter optimisation

For finding the best parameters for this model, the hyperparameters to tune PCA and the model were specified. Then a pipeline was created with 3 steps:

1. *scaler* which uses StandardScaler method to scale the data
2. *dimred* which applies PCA dimensionality reduction
3. *linearsvc* which applies the LinearSVC model.

Finally, a grid search was performed using the above-mentioned pipeline and the hyperparameter grid to find the combination with the best accuracy metric with the default cross-validation, that is, a 5-fold cross-validation.

```

1 from sklearn.svm import LinearSVC
2
3 param_grid = {
4     'dimred__n_components': [3, 4],
5     'linearsvc__C': [1e-2, 1e-1, 1, 10, 100],
6     'linearsvc__penalty': ['l1', 'l2'],
7     'linearsvc__dual': [False, True],
8     'linearsvc__class_weight': ['balanced', None]
9 }
10
11 linear_svc_pipeline = Pipeline([
12     ('scaler', scaler_pipeline),
13     ('dimred', PCA()),
14     ('linearsvc', LinearSVC())
15 ])
16
17 grid_search_lsvc = GridSearchCV(
18     linear_svc_pipeline,
19     param_grid=param_grid,
20     n_jobs=-1,
21     scoring='accuracy',
22     verbose=1
23 )
24
25 grid_search_lsvc.fit(X_train, y_train)

```

4.2.2 Model training

Then the model was trained using a similar pipeline, but this time with the optimal combination of parameters.

```

1 linear_svc_pipeline = Pipeline([
2     ('scaler', scaler_pipeline),
3     ('dimred', PCA(
4         n_components=grid_search_lsvc.best_params_['dimred__n_components']
5     )),
6     ('linearsvc', LinearSVC(
7         dual=grid_search_lsvc.best_params_['linearsvc__dual'],
8         C=grid_search_lsvc.best_params_['linearsvc__C'],
9         penalty=grid_search_lsvc.best_params_['linearsvc__penalty'],
10        class_weight=grid_search_lsvc.best_params_['linearsvc__class_weight']
11    ))
12 ])
13
14 linear_svc_pipeline.fit(X_train, y_train)

```

4.3 k-Nearest Neighbors Classifier

A k-Nearest Neighbors Classifier (KNeighborsClassifier) was implemented as part of this benchmark.

4.3.1 Hyperparameter optimisation

For finding the best parameters for this model, the hyperparameters to tune PCA and the model were specified. Then a pipeline was created with 3 steps:

1. *scaler* which uses StandardScaler method to scale the data
2. *dimred* which applies PCA dimensionality reduction
3. *kneighborsclassifier* which applies the KNeighborsClassifier model.

Finally, a grid search was performed using the above-mentioned pipeline and the hyperparameter grid to find the combination with the best accuracy metric with the default cross-validation, that is, a 5-fold cross-validation.

```

1 from sklearn.neighbors import KNeighborsClassifier
2
3 knn_pipeline = Pipeline([
4     ('scaler', scaler_pipeline),
5     ('dimred', PCA()),
6     ('kneighborsclassifier', KNeighborsClassifier())
7 ])
8
9 param_grid = {
10     'dimred__n_components': [3, 4],
11     'kneighborsclassifier__n_neighbors': range(3, 102, 3),
12 }
13
14 grid_search_knn = GridSearchCV(
15     knn_pipeline,
16     param_grid=param_grid,
17     n_jobs=-1,
18     scoring='accuracy',
19     verbose=1
20 )
21
22 grid_search_knn.fit(X_train, y_train)

```

4.3.2 Model training

Then the model was trained using a similar pipeline, but this time with the optimal combination of parameters.

```

1 knn_pipeline = Pipeline([
2     ('scaler', scaler_pipeline),
3     ('dimred', PCA(
4         n_components=grid_search_knn.best_params_['dimred__n_components']
5     )),
6     ('kneighborsclassifier', KNeighborsClassifier(
7         n_neighbors=grid_search_knn.best_params_['kneighborsclassifier__n_neighbors']
8     ))
9 ])
10
11
12 knn_pipeline.fit(X_train, y_train)

```

4.4 Artificial Neural Network

A Neural Network was also implemented as part of this benchmark. For this model, before optimising hyperparameters, it was necessary to define a function that created the proper neural network using Keras library.

4.4.1 Model definition

The amount of units and hidden layers in this model was specified during hyperparameter optimisation, as this function was defined with two arguments, namely *num_units* and *num_hidden_layers*, respectively defined with default values 10 and 1.

```

1 import tensorflow as tf
2 from keras.wrappers.scikit_learn import KerasClassifier
3 from keras.models import Sequential
4 from keras.layers import Dense
5
6 def build_clf(num_units=10, num_hidden_layers=1):
7     # initialising Sequential model and adding layers to it
8     ann_clf = tf.keras.models.Sequential()
9
10    # adding hidden layers
11    for i in range(num_hidden_layers):
12        ann_clf.add(tf.keras.layers.Dense(units=num_units, activation='relu'))
13
14    # adding output layer

```



```

15     ann_clf.add(tf.keras.layers.Dense(units=1, activation='sigmoid'))
16
17     # compiling model with chosen optimizer, loss function, and evaluation metrics
18     ann_clf.compile(
19         optimizer='adam',
20         loss='binary_crossentropy',
21         metrics=['accuracy']
22     )
23
24     return ann_clf

```

Also, the training data was split again to monitor the loss value and to find the best number of epochs on a further step, and this way this validation data was not considered by the model.

```

1 X_train_, X_val, y_train_, y_val = train_test_split(
2     X_train, y_train,
3     test_size=0.3, random_state=42
4 )
5
6 [X_train_.shape, X_val.shape, y_train_.shape, y_val.shape]

```

```
[(556518, 7), (238508, 7), (556518,), (238508,)]
```

4.4.2 Hyperparameter optimisation

For finding the best parameters for this model, the hyperparameters to tune the model were specified. Then a pipeline was created with 2 steps:

1. *scaler* which uses StandardScaler method to scale the data
2. *annmodel* which applies the Sequential model.

Finally, a grid search was performed using the above-mentioned pipeline and the hyperparameter grid to find the combination with the best accuracy metric with the default cross-validation, that is, a 5-fold cross-validation.

```

1 param_grid = {
2     'ann_model__num_units': [3, 4, 5, 6, 7],
3     'ann_model__num_hidden_layers': [2, 3, 4]
4 }
5
6 # creating an instance of the KerasClassifier using the defined function
7 ann_model = KerasClassifier(build_fn=build_clf, verbose=0)
8
9 ann_pipeline = Pipeline([
10     ('scaler', scaler_pipeline),
11     ('ann_model', ann_model)
12 ])
13
14 grid_search_ann = GridSearchCV(
15     ann_pipeline,
16     param_grid=param_grid,
17     n_jobs=-1,
18     scoring='accuracy',
19     verbose=0
20 )
21
22 grid_search_ann.fit(X_train_, y_train_)

```

In order to monitor the history of the loss evolution according the number of epochs, the *callback* parameter in the model *fit* method should be included, then the history could be assign to a variable. However, as the history object was not accessible by the pipeline, the model was scaled and trained separately to generate this history.

```

1 from keras import callbacks
2
3 earlystopping = callbacks.EarlyStopping(
4     monitor="val_loss", mode="min", patience=5, restore_best_weights=True,
5     verbose=-1

```

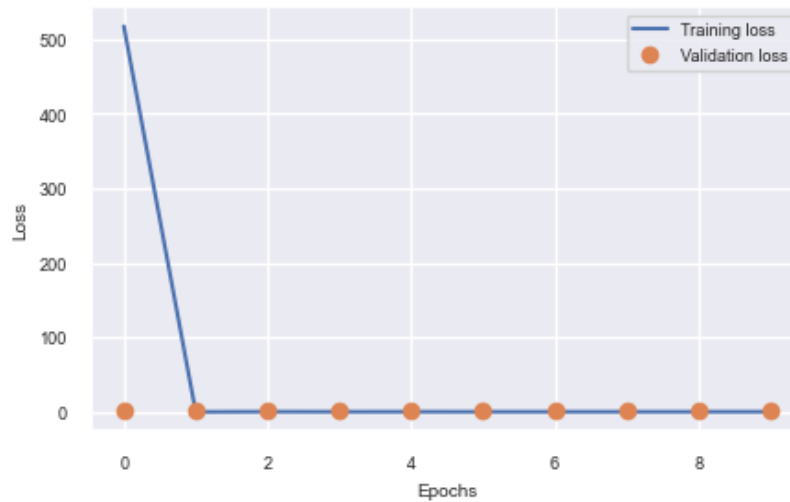


Figure 8: The training and validation loss for a neural network model trained over 15 epochs with early stopping

```

6 )
7
8 ann_model_ = KerasClassifier(
9     build_fn=build_clf,
10    num_hidden_layers=grid_search_ann.best_params_['ann_model__num_hidden_layers'],
11    num_units=grid_search_ann.best_params_['ann_model__num_units'],
12    verbose=1
13 )
14
15 scaler_pipeline.fit(X_val)
16 X_val = scaler_pipeline.transform(X_val)
17
18 history = ann_model_.fit(X_train_, y_train_,
19     epochs=15,
20     verbose=0,
21     validation_data=(X_val, y_val),
22     callbacks=[earlystopping]
23 )

```

The loss according the number of epochs on training and validation data can be seen on Figure 8 and a chart was created by the following code.

```

1 # plotting the training and validation loss
2 fig, (ax1) = plt.subplots(1, figsize=(5, 3))
3 ax1.plot(history.history['loss'])
4 ax1.plot(history.history['val_loss'], 'o')
5 ax1.set_xlabel('Epochs')
6 ax1.set_ylabel('Loss')
7 ax1.legend(['Training loss', 'Validation loss'])
8
9 plt.show()

```

4.4.3 Model training

Finally, the model was trained using a similar pipeline, with the optimal combination of parameters and the proper amount of epochs. The final configuration of this model is visible on Figure 9.

```

1 ann_model = KerasClassifier(
2     build_fn=build_clf,
3     num_units=grid_search_ann.best_params_['ann_model__num_units'],
4     num_hidden_layers=grid_search_ann.best_params_['ann_model__num_hidden_layers'],
5     verbose=0
6 )
7

```

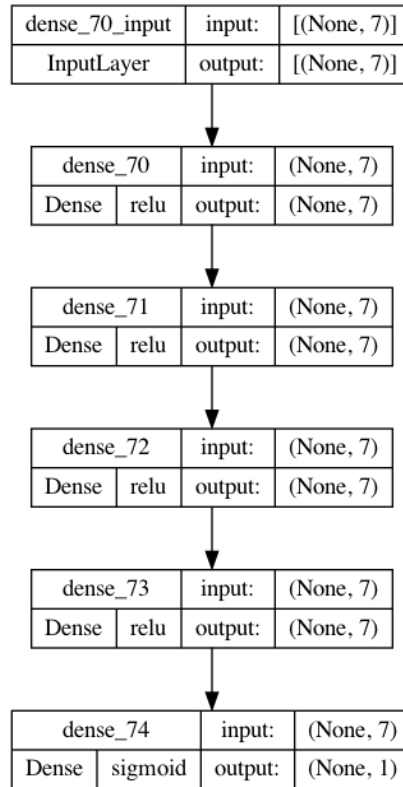


Figure 9: The architecture of a neural network model built using Keras

```

8 ann_pipeline = Pipeline([
9     ('scaler', scaler_pipeline),
10    ('ann_model', ann_model)
11 ])
12
13 ann_pipeline.fit(
14     X_train,
15     y_train,
16     ann_model__epochs=5,
17     ann_model__verbose=0
18 )

```

4.5 Decision Tree Classifier

As the data is not linearly separable and it's not possible to find a condition of 100% correctly classified by a hyperplane, a non-linear classifier is recommended, thus a Decision Tree Classifier was implemented as part of this benchmark.

4.5.1 Hyperparameter optimisation

For finding the best parameters for this model, the hyperparameters to tune the model were specified. Then a pipeline was created with 1 step:

1. *decisiontreeclassifier* which applies the Decision Tree model.

Finally, a grid search was performed using the above-mentioned pipeline and the hyperparameter grid to find the combination with the best accuracy metric with the default cross-validation, that is, a 5-fold cross-validation.

```

1 tree_pipeline = Pipeline([
2     ('decisiontreeclassifier', DecisionTreeClassifier())
3 ])

```

```

4
5 param_grid = {
6     'decisiontreeclassifier__min_samples_split': [2, 5, 10],
7     'decisiontreeclassifier__min_samples_leaf': [1, 2, 4],
8     'decisiontreeclassifier__max_features': ['sqrt', 'log2']
9 }
10
11 grid_search_tree = GridSearchCV(
12     tree_pipeline,
13     param_grid=param_grid,
14     n_jobs=-1,
15     scoring='accuracy',
16     verbose=1
17 )
18
19 grid_search_tree.fit(X_train, y_train)

```

4.5.2 Model training

Then the model was trained using a similar pipeline, but this time with the optimal combination of parameters.

```

1 tree_pipeline = Pipeline([
2     ('decisiontreeclassifier', DecisionTreeClassifier(
3         min_samples_split=grid_search_tree.best_params_['
4         decisiontreeclassifier__min_samples_split'],
5         min_samples_leaf=grid_search_tree.best_params_['
6         decisiontreeclassifier__min_samples_leaf'],
7         max_features=grid_search_tree.best_params_['
8         decisiontreeclassifier__max_features']
9     ))
10 ])
11 tree_pipeline.fit(X_train, y_train)

```

4.6 Random Forest Classifier

As the data is not linearly separable and it's not possible to find a condition of 100% correctly classified by a hyperplane, a non-linear classifier is recommended, thus a Random Forest Classifier was implemented as part of this benchmark.

4.6.1 Hyperparameter optimisation

For finding the best parameters for this model, the hyperparameters to tune the model were specified. Then a pipeline was created with 1 step:

1. *randomforestclassifier* which applies the Random Forest model.

Then, a grid search was performed using the above-mentioned pipeline and the hyperparameter grid to find the combination with the best accuracy metric with the default cross-validation, that is, a 5-fold cross-validation.

```

1 from sklearn.ensemble import RandomForestClassifier
2
3 rf_pipeline = Pipeline([
4     ('randomforestclassifier', RandomForestClassifier())
5 ])
6
7 param_grid = {
8     'randomforestclassifier__class_weight': ['balanced', 'balanced_subsample']
9 }
10
11 grid_search_rf = GridSearchCV(
12     rf_pipeline,
13     param_grid=param_grid,
14     n_jobs=-1,
15     scoring='accuracy',

```

```

16     verbose=1
17 )
18
19 grid_search_rf.fit(X_train, y_train)

```

4.6.2 Model training

Then the model was trained using a similar pipeline, but this time with the optimal combination of parameters.

```

1 rf_pipeline = Pipeline([
2     ('randomforestclassifier', RandomForestClassifier(
3         class_weight=grid_search_rf.best_params_['randomforestclassifier__class_weight']
4     ))
5 ])
6
7 rf_pipeline.fit(X_train, y_train)

```

5 Evaluation

5.1 Classification Report

The classification reports for all 5 models were generated by the code below.

```

1 # LinearSVC
2 cr_linesarvc = metrics.classification_report(y_test, y_predicted_lin_clf, digits=4)
3
4 # kNN
5 cr_knn = metrics.classification_report(y_test, y_predicted_knn, digits=4)
6
7 # Neural Networks
8 y_predicted_ann = y_predicted_ann.flatten()
9 y_predicted_ann = np.where(y_predicted_ann.round(2) > 0.5, 1, 0)
10 cr_ann = metrics.classification_report(y_test, y_predicted_ann, digits=4)
11
12 # Decision Tree
13 cr_tree = metrics.classification_report(y_test, y_predicted_tree, digits=4)
14
15 # Random Forest
16 cr_rf = metrics.classification_report(y_test, y_predicted_rf, digits=4)

```

5.1.1 LinearSVC

```

1 # printing classification report for LinearSVC
2 print(cr_linesarvc)

```

	precision	recall	f1-score	support
0	0.9980	0.9808	0.9893	2763432
1	0.7568	0.9687	0.8498	170839
accuracy			0.9801	2934271
macro avg	0.8774	0.9747	0.9195	2934271
weighted avg	0.9840	0.9801	0.9812	2934271

5.1.2 k-Neighbors Classifier

```

1 # printing classification report for kNN classifier
2 print(cr_knn)

```

	precision	recall	f1-score	support
0	0.9999	0.9949	0.9974	2763432
1	0.9233	0.9979	0.9591	170839
accuracy			0.9950	2934271
macro avg	0.9616	0.9964	0.9782	2934271
weighted avg	0.9954	0.9950	0.9951	2934271

5.1.3 Neural Networks

```
1 # printing classification report for ANN
2 print(cr_ann)
```

	precision	recall	f1-score	support
0	0.9980	0.9831	0.9905	2763432
1	0.7795	0.9688	0.8639	170839
accuracy			0.9822	2934271
macro avg	0.8888	0.9759	0.9272	2934271
weighted avg	0.9853	0.9822	0.9831	2934271

5.1.4 Decision Tree

```
1 # printing classification report for Decision Tree
2 print(cr_tree)
```

	precision	recall	f1-score	support
0	1.0000	0.9998	0.9999	2763432
1	0.9960	0.9998	0.9979	170839
accuracy			0.9998	2934271
macro avg	0.9980	0.9998	0.9989	2934271
weighted avg	0.9998	0.9998	0.9998	2934271

5.1.5 Random Forest

```
1 # printing classification report for Random Forest
2 print(cr_rf)
```

	precision	recall	f1-score	support
0	1.0000	0.9999	1.0000	2763432
1	0.9987	1.0000	0.9993	170839
accuracy			0.9999	2934271
macro avg	0.9993	1.0000	0.9996	2934271
weighted avg	0.9999	0.9999	0.9999	2934271

5.2 Confusion Matrices

In this section the confusion matrices for every model for test labels were computed and plotted side by side on Figure 10.

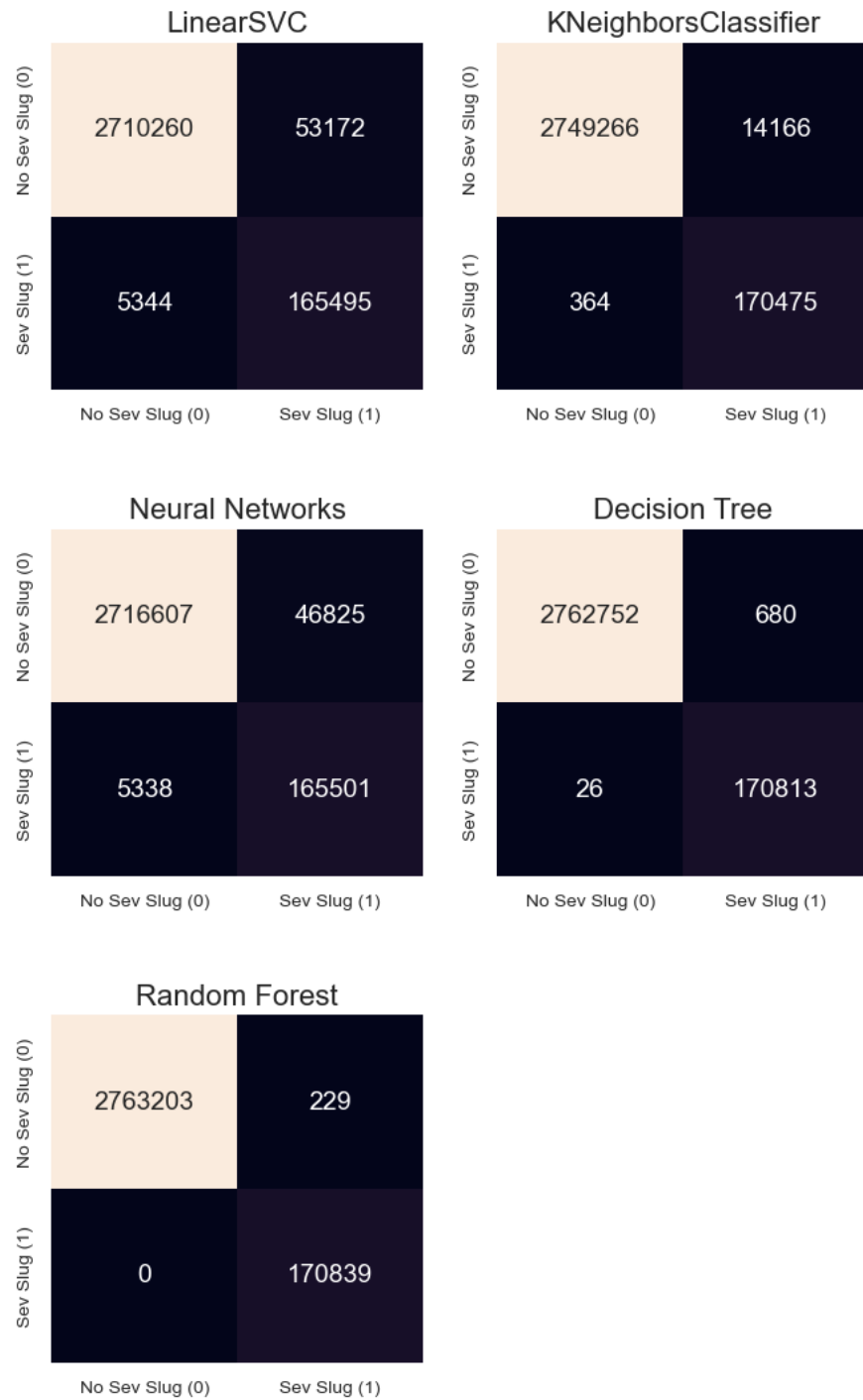


Figure 10: Confusion matrices showing the performance of five machine learning models in predicting severe slugging occurrences

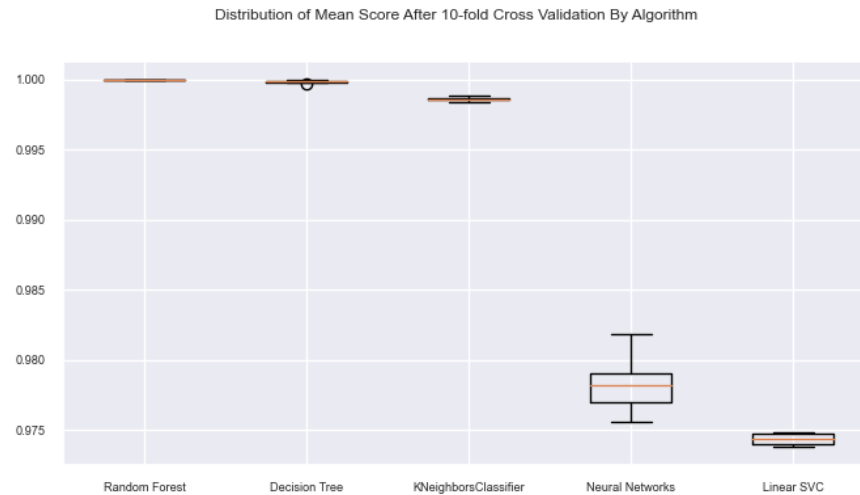


Figure 11: Box plots showing the distribution of Mean Score After 10-fold Cross Validation By Algorithm

5.3 10-Fold Cross Validation

Finally, a K-Folds cross-validator was selected to split training data set into 10 consecutive folds and compute the mean accuracy and the respective standard deviation for each model. The models were still put on the same pipelines used on modeling section.

```
1 models = [
2     ('Random Forest', rf_pipeline),
3     ('Decision Tree', tree_pipeline),
4     ('KNeighborsClassifier', knn_pipeline),
5     ('Neural Networks', ann_pipeline),
6     ('Linear SVC', linear_svc_pipeline)
7 ]

1 for name, model in models:
2     kfold = KFold(n_splits=10, shuffle=True, random_state=42)
3
4     cv_results = cross_val_score(model, X_train, y_train, cv=kfold, scoring='accuracy'
5     )
6     cross_validation_dict[name] = cv_results
7
8     msg = "%s: %f (%f)" % (name, np.nanmean(cv_results), np.nanstd(cv_results))
9     print(msg)
```

```
Random Forest: 0.999940 (0.000018)
Decision Tree: 0.999826 (0.000069)
KNeighborsClassifier: 0.998610 (0.000138)
Neural Networks: 0.978187 (0.001739)
Linear SVC: 0.974343 (0.000403)
```

The following code was implemented to plot how the mean score after cross validation was distributed in each model. This chart can be visualised on Figure 11.

```
1 # comparing algorithms
2 fig = plt.figure()
3 fig.suptitle('Distribution of Mean Score After 10-fold Cross Validation By Algorithm')
4 ax = fig.add_subplot(111)
5 plt.boxplot(cross_validation_dict.values())
6 ax.set_xticklabels(cross_validation_dict.keys())
7 fig.set_size_inches(8,4)
8 plt.show()
```


6 Conclusion

In this study, Random Forest and Decision Tree Classifiers had the best performance overall with 99.994% and 99.982% of accuracy respectively. These non-linear classifiers could find better results because the 3W data set does not present a clear linear separability between records related to Severe Slugging and other normal records or other undesirable events.

References

Vargas, Ricardo Emanuel Vaz et al. (July 2019). *A realistic and public dataset with rare undesirable real events in oil wells*. URL: <https://www.sciencedirect.com/science/article/pii/S0920410519306357?via=ihub>.