
Assignment - 2 (Weather Classification)

Team Members:

Achintya Harsha (IMT2021525)

Adithya Nagaraja Somasalle (IMT2021054)

Devendra Rishi Nelapati (IMT2021076)

December 17, 2023



Contents

1	Exploratory Data Analysis (EDA)	2
1.1	Dealing with Null Values	2
1.2	Distribution of Features	3
1.3	Correlation	6
1.4	Outlier Detection and removal	7
2	Pre-processing of Data	8
2.1	General Data Preprocessing	8
2.2	Location-Based Priority Probability Calculation	8
2.3	Seasonal Data Categorization	8
2.4	Upsampling and Downsampling	8
2.5	Categorical Data Encoding	9
2.6	Dealing with Periodic Features	9
3	Model training and testing	9
3.1	SVM Classifier	9
3.2	Neural Networks Classifier	10
3.3	Ensemble XGBoost Classifier (Best Model)	10
4	Timeline	11

1 Exploratory Data Analysis (EDA)

1.1 Dealing with Null Values

First we look for null values in the dataframe. Using the library missingno, we can visualise the percentage of missing values in each column.

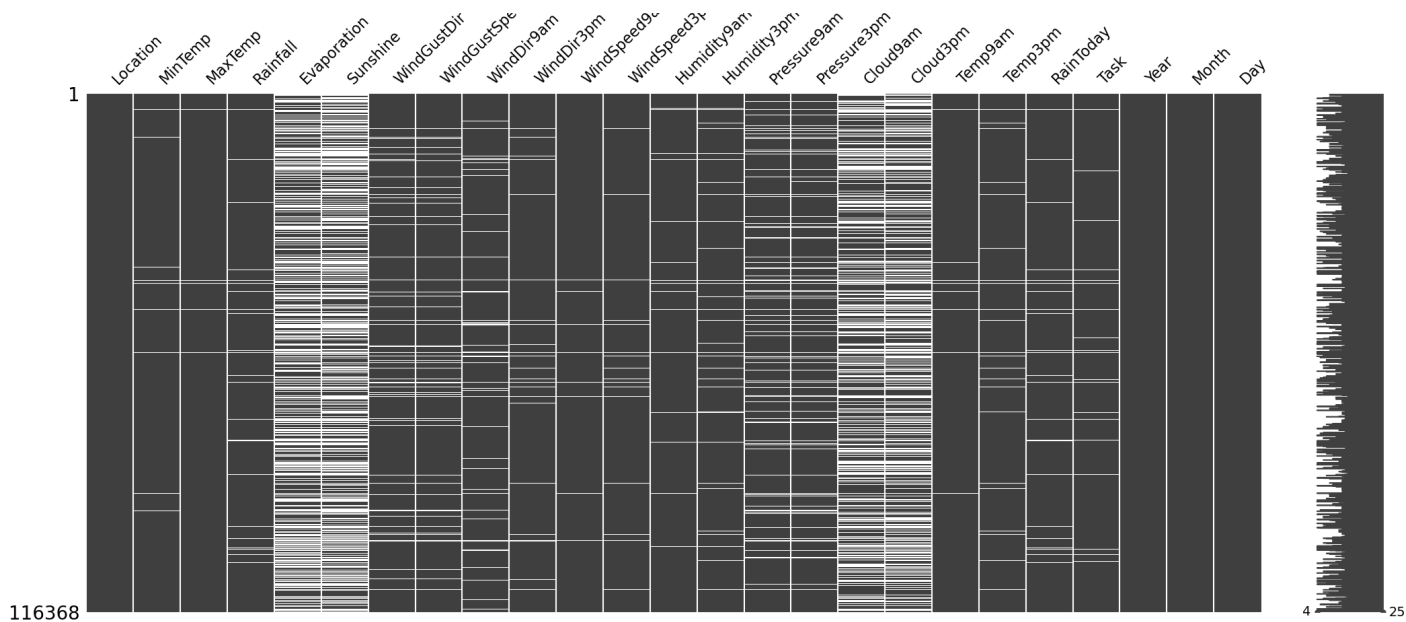


Figure 1: *Missing values in Dataframe*

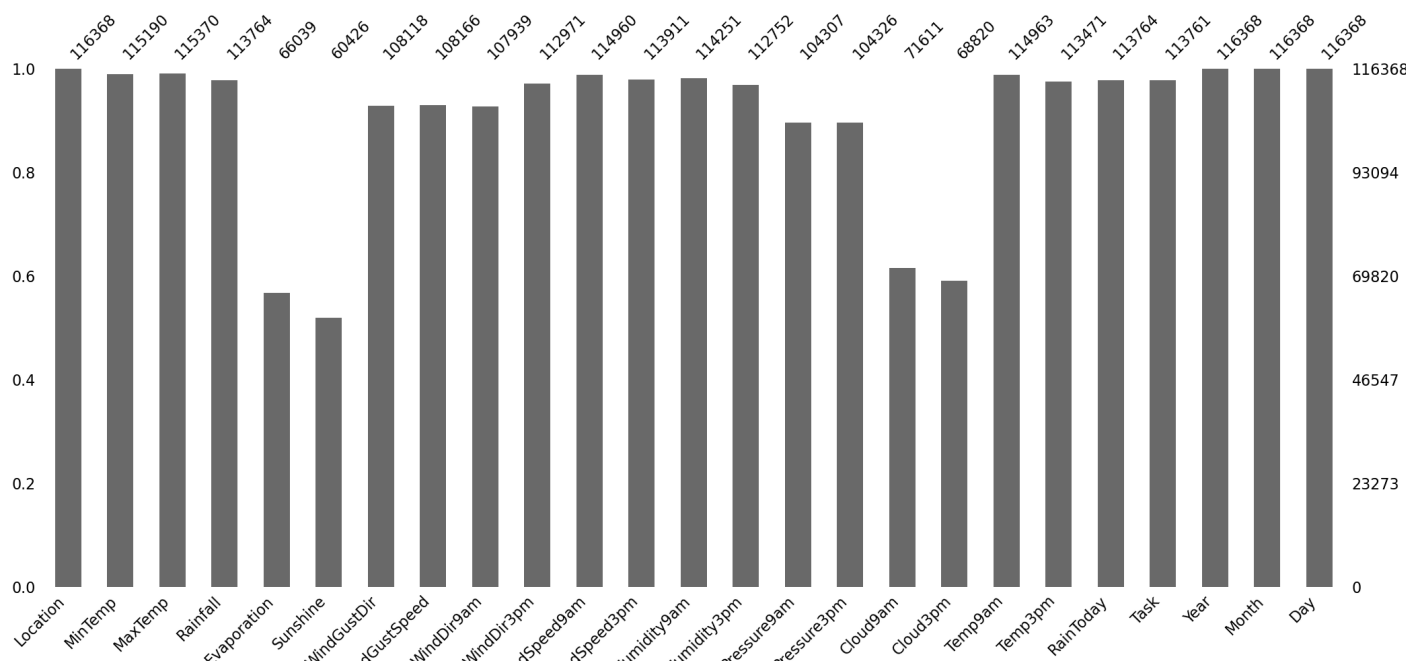


Figure 2: *Percentage of Missing values in Dataframe*

From the graphs we observe that columns Sunshine, Evaporation, Cloud9am, Cloud3pm has large percent of null values. Among these, Sunshine and Evaporation have null percentage nearly equal to

50%.

So we drop the respective columns from the data-set such a high percent of null values can cause inaccuracy in the prediction model.

1.2 Distribution of Features

Here we will probe into the distribution of the values of different features. We will look at the mean and variance of a particular feature and the probability of the occurrence of a particular value.

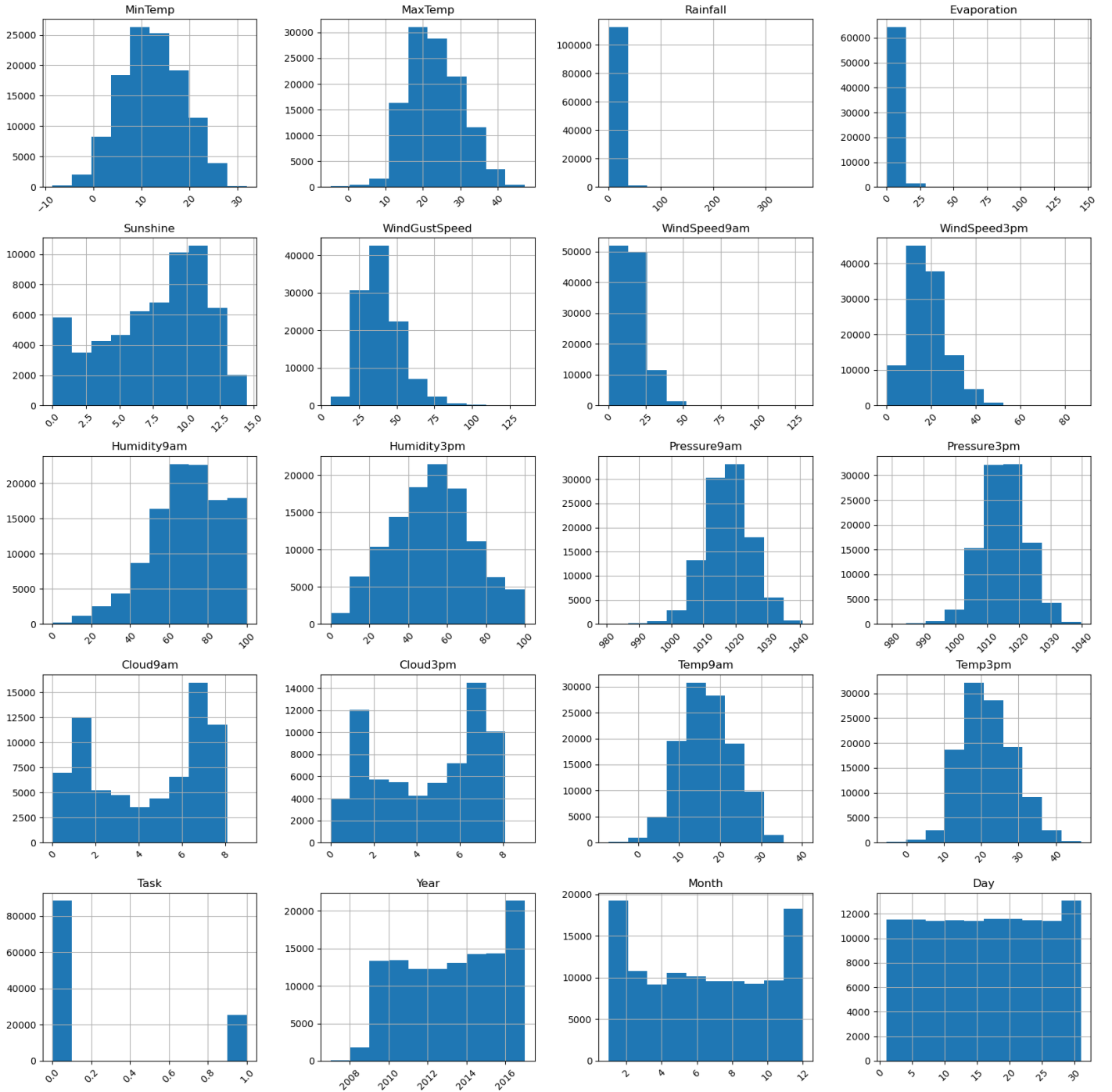


Figure 3: Histogram showing the mean and spread of numerical features

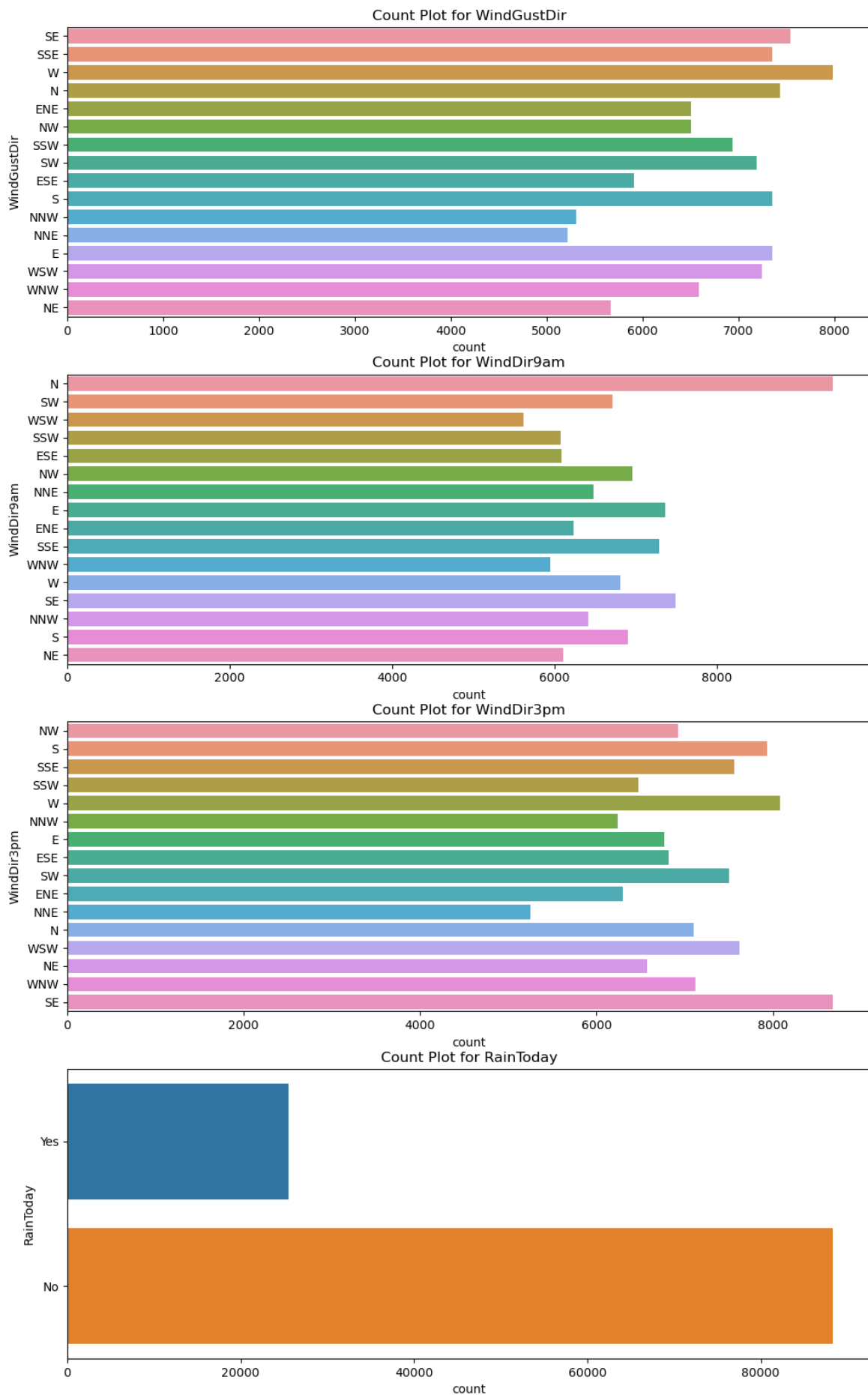


Figure 4: *Count plot showing the mean and spread of categorical features*

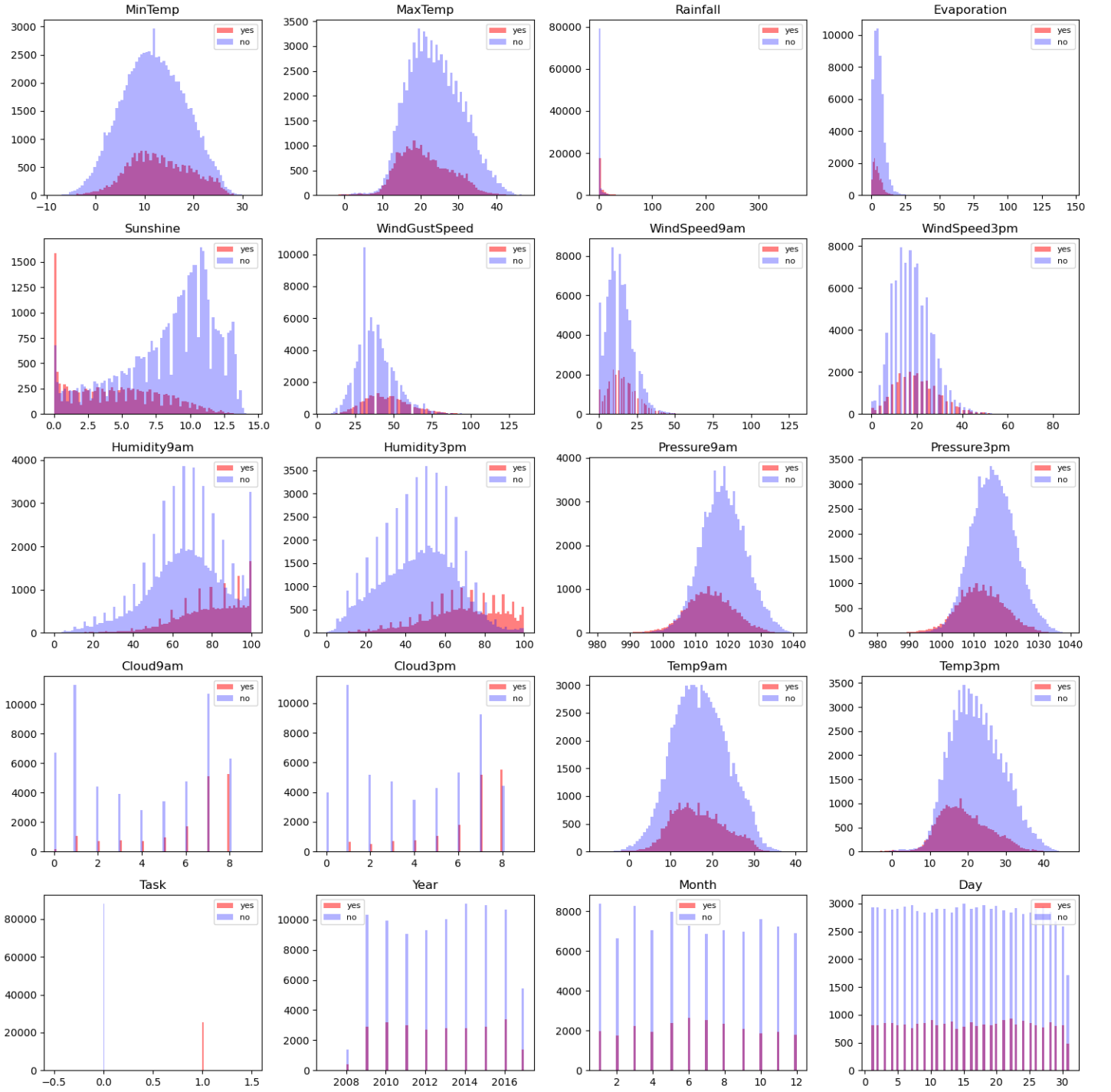


Figure 5: *Comparison of Characteristics for Task Values 0 and 1*

From the above graphs we can remove the features that have low variance. Also if we look that Figure5, some features like 'MinTemp', 'MaxTemp', 'Day' have similar mean and variance for cases where Task is 0 and 1. So they are not really useful for prediction.

1.3 Correlation

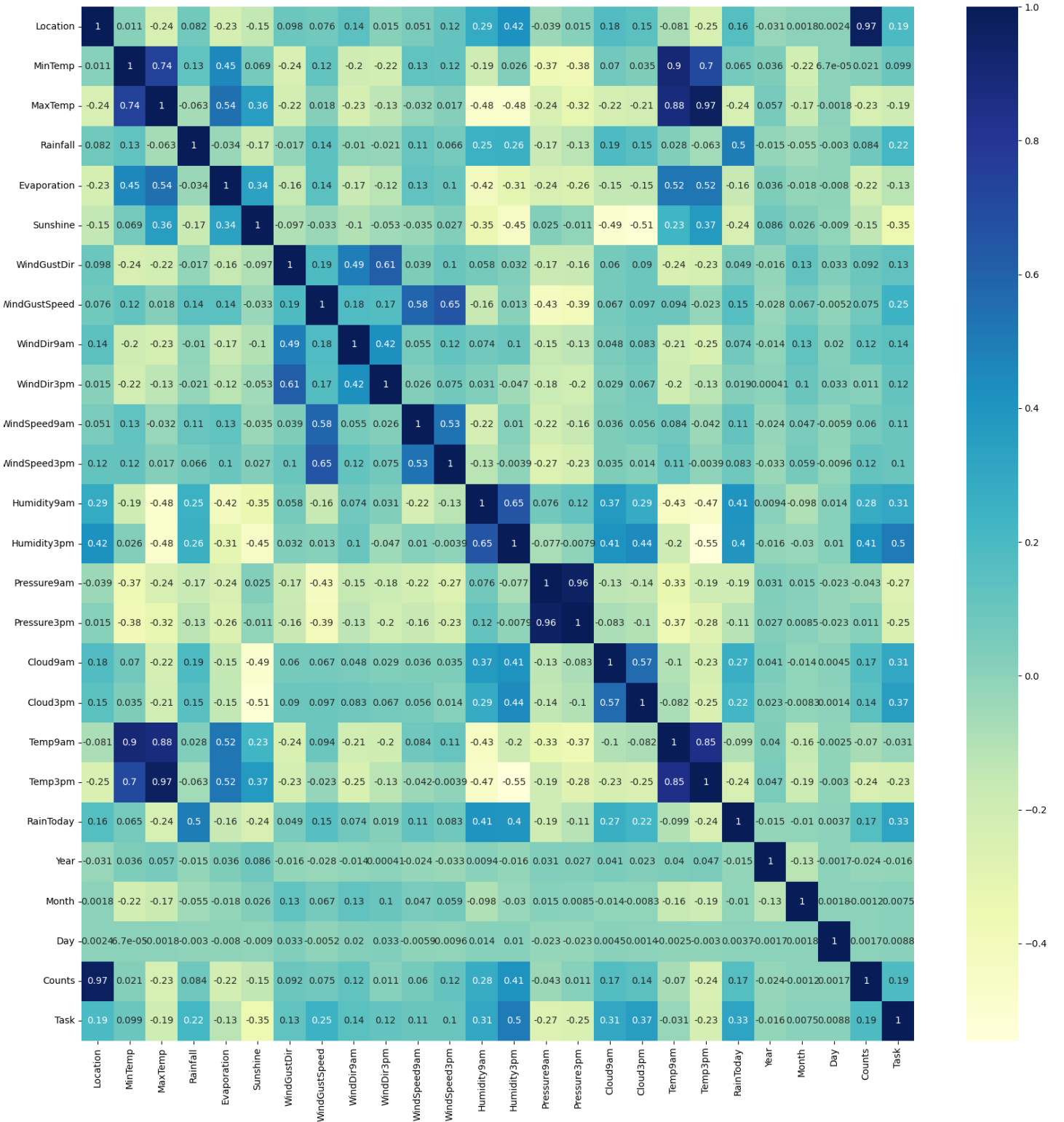


Figure 6: Correlation heat map of every feature w.r.t. each other

Features that are poorly correlated with Task(correlation close to 0) are removed as they do not contribute to improving accuracy. Example: 'Day' feature is removed.

1.4 Outlier Detection and removal

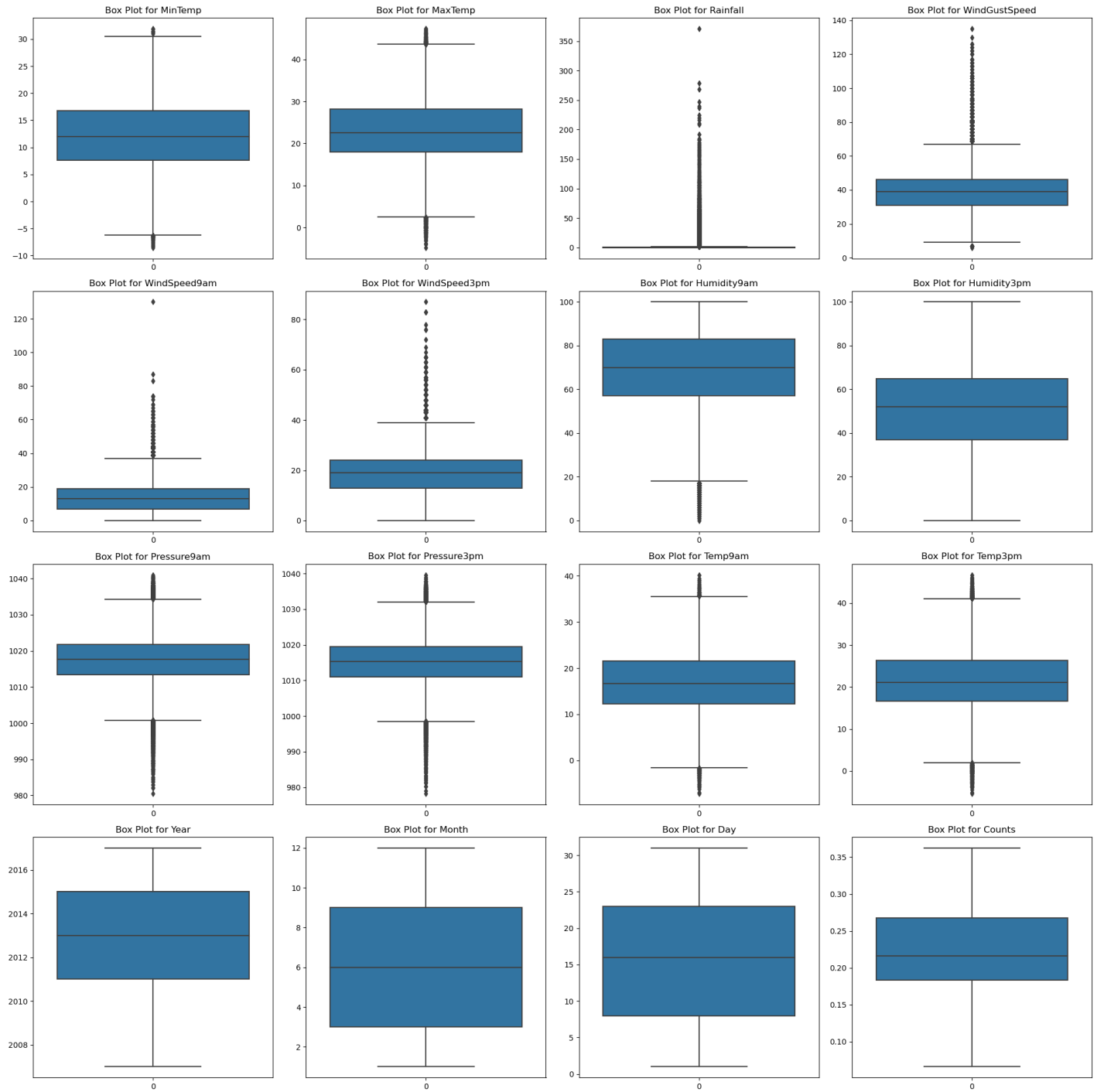


Figure 7: *Outliers plotted using boxplot*

In our analysis, we identified and addressed outliers in the dataset. We either removed them to improve data quality and, potentially, enhance model performance.

2 Pre-processing of Data

2.1 General Data Preprocessing

We began by examining the provided dataset, consisting of test and train CSV files. One of our initial steps involved breaking down the date column into more categorical columns, such as date, month, and year. This action was taken to manage the potential issue of the date column having a large number of unique values. Consequently, we simplified it into 31 days, 12 months, and several years.

```
len(pd.unique(train_df['Date']))  
3414
```

Figure 8: *Number of unique date values without splitting is very large*

We also removed rows with null values in columns where the null value percentage was less than 5%. These rows were considered inconsequential as the columns contained sufficient data for training our model.

2.2 Location-Based Priority Probability Calculation

Subsequently, we conducted priori probability calculations based on location. This involved dividing the count of each location value by its location coin value when the Task value of that particular row equaled 1. This calculation provided us with the probability that the task value would be 1 at a given location. Also total number of occurrences of a particular location was calculated. Using this we calculated the probability of target being 1 given the place. The resulting probability values were then added as a new column(Counts) to the data set.

$$trainDF[Counts] = P(Task = 1|Location)$$

2.3 Seasonal Data Categorization

Given that the task in question is 'Humidity,' it's a logical assumption that humidity is influenced by the season. Thus, we organized the data frame into four sub-dataframes, categorized by seasons: Spring_df (March, April, May), Summer_df (June, July, August), Fall_df (September, October, November), and Winter_df (December, January, February).

For data imputation, we addressed the missing values within the respective sub-dataframes. For numerical data, we replaced missing values with the mean of the column values for that particular sub-data frame, while for categorical data, we used the mode of the sub-data frame.

2.4 Upsampling and Downsampling

Another noteworthy observation was the substantial disparity in the number of rows with task value 0 compared to those with task value 1. To address this, we employed a combination of up-sampling and

down-sampling. We divided the dataframe into two sub-dataframes based on the task value, then found the midpoint of the total rows. We up-sampled the task=1 rows to match this value and down-sampled the task=0 rows accordingly. Subsequently, we merged these sub-dataframes into a single one, ensuring the rows were shuffled to maintain randomness.

2.5 Categorical Data Encoding

Further, we had to deal with categorical data that required encoding. Initially, we attempted one-hot encoding, which resulted in an accuracy of approximately 20%. However, after implementing label encoding, the accuracy score improved to 64%. Further enhancements were achieved by employing target encoding, where each unique value in a categorical column was assigned a unique float value. This approach elevated the accuracy to 68%.

2.6 Dealing with Periodic Features

We noticed that the Data, Month and Wind Direction are Periodic features. So to introduce periodicity in the data so that the model can interpret it, we had to apply a cos function on the features to make them cyclic.

It's crucial to note that all these preprocessing steps were also applied to the test data. The primary difference was that for the test data, we performed null value imputation directly on the dataset to prevent any shuffling of rows.

3 Model training and testing

3.1 SVM Classifier

SVM Classifier implementation code

```
1 from sklearn.svm import SVC
2 from sklearn.metrics import accuracy_score, f1_score
3 from sklearn.model_selection import train_test_split
4
5 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.1, random_state=42)
6
7 svm_model = SVC(kernel='rbf', C=1, gamma='scale')
8
9 svm_model.fit(X_train, y_train)
10
11 y_pred = svm_model.predict(X_test)
12
13 accuracy = f1_score(y_test, y_pred)
14 print(f'Accuracy: {accuracy * 100:.2f}%')
```

Best Accuracy : 62.211%

3.2 Neural Networks Classifier

Neural Networks Classifier implementation code

```
1
2 model = keras.Sequential([
3     keras.layers.Dense(256, activation='relu', input_shape=(X_train_preprocessed.shape[1],)),
4     keras.layers.Dropout(0.2),
5     keras.layers.Dense(192, activation='relu'),
6     keras.layers.Dropout(0.2),
7     keras.layers.Dense(256, activation='relu'),
8     keras.layers.Dropout(0.2),
9     keras.layers.Dense(32, activation='relu'),
10    keras.layers.Dropout(0.2),
11    keras.layers.Dense(160, activation='relu'),
12    keras.layers.Dropout(0.2),
13    keras.layers.Dense(1, activation='sigmoid')
14 ])
15
16 model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
17
18 early_stopping = EarlyStopping(monitor='val_loss', patience=10, restore_best_weights=True)
19
20 history = model.fit(
21     X_train_preprocessed,
22     y_train,
23     epochs=150,
24     batch_size=64,
25     validation_split=0.2,
26     callbacks=[early_stopping],
27     verbose=2
28 )
29
30 test_loss, test_acc = model.evaluate(X_test_preprocessed, y_test)
31 print(f'Test accuracy: {test_acc}')
```

Best Accuracy : 67.160%

3.3 Ensemble XGBoost Classifier (Best Model)

XGBoost Classifier implementation code

```
1 import xgboost as xgb
2
3 models = []
4 np.random.seed(1)
5 for i in range(10):
6     model = xgb.XGBClassifier(
7         learning_rate=0.05,
8         n_estimators=750 + (0 if (i == 0) else np.random.randint(-20, 20)),
9         max_depth=8,
10        min_child_weight=1,
11        gamma=0.5,
12        subsample=0.6,
13        colsample_bytree=1,
```

```

14         reg_alpha=0.1,
15         random_state=42,
16         reg_lambda=2
17     )
18     model.fit(X_train, y_train)
19     models.append(model)
20
21     ## Prediction
22     predictions = [model.predict(test_df) for model in models]
23     Y_pred = np.zeros(test_df.shape[0])
24
25     for i in range(len(models)):
26         for j in range(test_df.shape[0]):
27             Y_pred[j] += predictions[i][j]
28
29     Y_pred = Y_pred > len(models)/2

```

Best Accuracy : 69.353%

4 Timeline

XGBoost turned out to be the best model for the given dataset.

- **62.211%** : SVM 3 was implemented with hyper-parameter $C = 1$ and kernel as Radial Basis Function.
- **65.645%** : ANN model was built and trained on the data. One hot encoding was used and the columns were scaled. We noticed that there was an imbalance in the the number of Task = 0 and 1. So we implemented up and down sampling.
- **67.160%** : We noticed that the Data, Month and Wind Direction are Periodic features. So to introduce periodicity in the data so that he model understands it, we had to apply a cos function on the features to make them cyclic.
- **68.518%** : We switched to XGBoost for better classification. Seasonal Data Categorization was implemented
- **69.353%** : Ensemble learning was used along with the XGBoost built previously to achieve this accuracy.