

data_cleaner_with_test_df_functionality.DataCleaner

```
class data_cleaner_with_test_df_functionality.DataCleaner(dataframe, features,  
test_df=None, mode="wnv")
```

Cleans the data provided in GA DSIF Project 4.

This class takes in a dataframe and a chosen feature list created from merged csv files from the GA DSIF project 4 assets, and creates an object with the attributes df, features, X, y, X_train, y_train, X_test, y_test.

If the use of a separate test df is intended, test df may be passed as an optional keyword argument.

Most attributes are initialised as None, and are filled by executing the clean() method on the object. The clean method contains a series of functions which start by removing outliers and populating the X and y attributes of the object based on the list of features passed in.

Each function thereafter is programmed to handle exceptions where the column has been removed due to it not being in X as a result of the features passed in. Please refer to the comments relating to each function below for a further explanation of what each function does.

The clean_test method may only be executed after the clean method has been executed as the encoder and scalers must first be fitted on training data.

The class object can be instantiated in two modes. In "num_mosquitos" mode, it is intended to clean (1) the train dataset with "NumMosquitos" as the target variable, and (2) the test dataset with the intention of predicting "NumMosquitos". In "wnv" mode, it is intended to clean (1) the train dataset with "wnvpresent" as the target variable, and (2) the test dataset with the intention of predicting "wnvpresent". This mode function was developed due to "NumMosquitos" being absent from the provided test dataset. The DataCleaner object will default to "wnv" mode if no argument is passed.

Parameters:	dataframe: training dataframe
	Training dataset in the form of a pandas dataframe.
	features: list of features
	Features to be processed in the form of a python list.

	<p>test_df: testing dataframe, optional</p> <p>Testing dataset in the form of a pandas dataframe.</p> <p>mode: {'wnv', 'num_mosquitos'}, default = 'num_mosquitos'</p> <p>Sets the mode which the class constructor will operate in.</p>
Attributes:	<p>df: dataframe</p> <p>Training dataframe passed at instantiation.</p> <p>test_df: test_df</p> <p>Optional testing dataframe passed at instantiation.</p> <p>features: python list</p> <p>List of features to be processed passed at instantiation.</p> <p>categorical: python list</p> <p>List of categorical features found in the datasets: ["month", "Species", "ResultDir", "gps_cat", "sprayed", "Sunrise", "Sunset"].</p> <p>cat_features: python list</p> <p>List of categorical features to be processed, dependent upon features.</p> <p>ohe: object</p> <p>sk-learn OneHotEncoder instance.</p> <p>encoded_features: python list</p> <p>List of features after one hot encoding.</p> <p>encoded_test_features: python list</p> <p>List of features of test_df after one hot encoding.</p> <p>mode: str</p>

	<p>Toggle selector to operate the class in either wnv mode or num_mosquitos mode.</p> <p>y: pandas series</p> <p>Training target variable.</p> <p>X: pandas dataframe</p> <p>Training predictors.</p> <p>X_train: pandas dataframe</p> <p>Training predictors after application of sk-learn's train-test-split.</p> <p>X_test: pandas dataframe</p> <p>Testing predictors derived from training dataset after application of sk-learn's train-test-split.</p> <p>y_train: pandas dataframe</p> <p>Training target after application of sk-learn's train-test-split.</p> <p>y_test: pandas dataframe</p> <p>Testing target derived from training dataset after application of sk-learn's train-test-split.</p> <p>clean_executed: boolean</p> <p>Flags whether the .clean() method has been executed.</p> <p>num_mosquito_ss: object</p> <p>sk-learn standard scaler object fitted to the relevant column in the training dataset, saved for transforming relevant testing data column.</p> <p>tmax_ss: object</p> <p>sk-learn standard scaler object fitted to the relevant column in the training dataset, saved for transforming relevant testing data column.</p> <p>self.tmin_ss: object</p>
--	--

	<p>sk-learn standard scaler object fitted to the relevant column in the training dataset, saved for transforming relevant testing data column.</p> <p>self.tavg_ss: object</p> <p>sk-learn standard scaler object fitted to the relevant column in the training dataset, saved for transforming relevant testing data column.</p> <p>self.dewpoint_ss: object</p> <p>sk-learn standard scaler object fitted to the relevant column in the training dataset, saved for transforming relevant testing data column.</p> <p>self.wetbulb_ss: object</p> <p>sk-learn standard scaler object fitted to the relevant column in the training dataset, saved for transforming relevant testing data column.</p> <p>self.heat_ss: object</p> <p>sk-learn standard scaler object fitted to the relevant column in the training dataset, saved for transforming relevant testing data column.</p> <p>self.cool_ss: object</p> <p>sk-learn standard scaler object fitted to the relevant column in the training dataset, saved for transforming relevant testing data column.</p> <p>self.preciptotal_ss: object</p> <p>sk-learn standard scaler object fitted to the relevant column in the training dataset, saved for transforming relevant testing data column.</p> <p>self.windspeed_ss: object</p> <p>sk-learn standard scaler object fitted to the relevant column in the training dataset, saved for transforming relevant testing data column.</p> <p>self.pressure_ss: object</p> <p>sk-learn standard scaler object fitted to the relevant column in the training dataset, saved for transforming relevant testing data column.</p> <p>self.sealevel_ss: object</p>
--	--

	sk-learn standard scaler object fitted to the relevant column in the training dataset, saved for transforming relevant testing data column.
--	---

Examples:

```
import pandas as pd
from data_cleaner_with_test_df_functionality import DataCleaner

features = ['Species', 'Tmax', 'Tmin', 'Tavg', 'DewPoint', 'WetBulb',
            'Heat', 'Cool', 'PrecipTotal',
            'ResultSpeed', 'ResultDir', 'StnPressure', 'SeaLevel',
            'gps_cat', 'sprayed', 'month']

df = pd.read_csv("df.csv", low_memory=False)
test_df = pd.read_csv("test_df.csv", low_memory=False)

cleaned = DataCleaner(df, features, test_df, mode="wnv")

cleaned.clean()
```

```
import pandas as pd
from data_cleaner_with_test_df_functionality import DataCleaner

features = ['Species', 'Tmax', 'Tmin', 'Tavg', 'DewPoint', 'WetBulb',
            'Heat', 'Cool', 'PrecipTotal',
            'ResultSpeed', 'ResultDir', 'StnPressure', 'SeaLevel',
            'gps_cat', 'sprayed', 'month']

df = pd.read_csv("df.csv", low_memory=False)

cleaned = DataCleaner(df, features)

cleaned.clean()
```

Methods

.clean()	Cleans training dataframe.
.clean_test()	Cleans testing dataframe. Must have executed .clean() prior to executing this method.