

Anomaly Detection Documentation

windows.py

class Window

Object	object object
Methods	None; to be updated
Attributes	None; to be updated

class DataFrameWin

Object	Window object
Methods	<p><code>__init__</code>:</p> <ul style="list-style-type: none"><i>df</i> or <i>csv_file</i>: only 1 is optional, if both are provided then <i>df</i> is used<i>df</i>: Pandas DataFrame<i>csv_file</i>: Directory to csv file to use as dataframe<i>drop_columns</i>: optional list, only compatible with <i>csv_file</i>, drops columns in list<i>parse_dates</i>: optional string, name of column of timestamps to parse <p>modify:</p> <ul style="list-style-type: none"><i>mode</i>: string, either “drop” or “format_dates”<ul style="list-style-type: none"><i>drop</i>: <i>value</i> parameter can be either string or list of strings of column name(s) to drop<i>format_dates</i>: <i>value</i> parameter should be string, list of timestamp column to format<ul style="list-style-type: none">timestamp column should be any comprehensible time formatthe function modifies by splitting the column into three other columns for time/day/monthtime: time past midnight, day: day of week from Monday (0) to Sunday (6), month: (1-12)<i>value</i>: this parameter’s expectation depends on <i>mode</i> (see above) <p>create_lag:</p> <ul style="list-style-type: none"><i>n</i>: number of lags to create; a lag is simply a column shifted downwards, each lag will be shifted incrementally<i>col</i>: the column to create the lags with; this will be left intact <p>interpolate_missing_rows:</p> <ul style="list-style-type: none">No parameters, interpolates missing dataframe rows <p>split_by_col:</p> <ul style="list-style-type: none"><i>col_name</i>: The name of the column to split by (DataFrameWin objects for each different value in the column)Returns: Dictionary; keys are the unique values in <i>col_name</i>; values are DataFrameWin objects
Attributes	<p><i>self.dates_formatted</i>: True is dates are formatted using <i>modify</i>(“format_dates”,...) and false otherwise</p> <p><i>self.show_all</i>: If set to true, the next time (and only next time) you print the object, it will show all the rows/columns</p> <p><i>self.nlags</i>: Keeps track of number of lags in function, DO NOT change this</p>

class GraphWin

Object	Window object
Methods	<p><code>__init__</code>:</p> <ul style="list-style-type: none">DataFrameWin: DataFrameWin object to connect graphing window to <p>export_DF:</p> <ul style="list-style-type: none">Returns: DataFrameWin object the graph window is connected to <p>create_graph:</p> <ul style="list-style-type: none"><i>mode</i>: string, either “missing-values” or “pacf” or “heatmap”<ul style="list-style-type: none"><i>missing_values</i>: Creates missing value histogram with 50 bins<ul style="list-style-type: none"><i>x</i> should be the x axis<i>y</i> should be the frequency column (how many missing values in y for each x?)<i>pacf</i>: Creates partial autocorrelation plot<ul style="list-style-type: none"><i>x</i> should be the x axis

	<ul style="list-style-type: none"> ▪ y should be an integer in this case; the x-scale ○ <i>heatmap</i>: Creates a heatmap with time of day/day of week; DataFrameWin MUST be time formatted <ul style="list-style-type: none"> ▪ x should be the value column for heatmap values ▪ y is not used, set to None • x: depends on mode, see above • y: depends on mode, see above
Attributes	self.DataFrameWin: DataFrameWin object graphing window is connected to

class TrainingWin

Object	Window object
Methods	<p><u>__init__</u></p> <ul style="list-style-type: none"> • DataFrameWin: DataFrameWin object to connect training window to <p>export_DF:</p> <ul style="list-style-type: none"> • Returns: DataFrameWin object the graph window is connected to <p>split:</p> <ul style="list-style-type: none"> • x_attr: The X attributes (feature columns) • y_attr: The y attributes (label columns) • Creates X/y train/test split based on the above (saved within object) <p>tune_hps:</p> <ul style="list-style-type: none"> • param_grid: A dictionary with the keys being the possible tuning parameters and keys being dictionary with values to consider (hyper parameters saved within object) <ul style="list-style-type: none"> ○ Possible parameters (keys): 'bootstrap', 'max_depth', 'max_features', 'min_samples_leaf', 'min_samples_split', 'n_estimators' <p>generate_model:</p> <ul style="list-style-type: none"> • Creates model based on saved hyper_parameters, if no hyper parameters are saved then use default • Model saved within object <p>save_model:</p> <ul style="list-style-type: none"> • dir_: directory to save model as “.sav” file in; “.sav” must be included within the directory <p>load_model:</p> <ul style="list-style-type: none"> • dir_: file directory to load “.sav” file from; model stored within object <p>predict_csv:</p> <ul style="list-style-type: none"> • Predicts csv file using model and exports as a separate csv • dir1: Directory of first csv file (one to predict out of) • timeCol: String; name of timestamp column • valueCol: String; name of value column • dir2: Directory of second csv file (one to export to) • anomaly: bool or int; optional, if set as int then create an anomaly column, if set as False then don't <ul style="list-style-type: none"> ○ The int represents how strict anomaly detection is (larger = less rows are categorized as anomalies) <ul style="list-style-type: none"> ▪ $\text{forecasted} - \text{ground truth} > \text{anomaly} * \text{Standard Deviation} \implies \text{Anomaly}$ <p>predict_list:</p> <ul style="list-style-type: none"> • Returns prediction given list of features • lis: The list of feature columns' values (must be same order given in split method) • Returns: the prediction given the feature columns' values
Attributes	<ul style="list-style-type: none"> • self.DataFrameWin: DataFrameWin object graphing window is connected to • self.X_train, self.X_test, self.y_train, self.y_test: The train/test/x/y splits • self.hyper_param: The optimal hyper parameters • self.model: The model that drives everything

Example Code

```
from windows import Window
from windows import DataFrameWin
from windows import GraphWin
from windows import TrainingWin

df = DataFrameWin(csv_file=r"C:\path\to\file.csv") # create dataframe from CSV

df.modify("format_dates", "timestamp") # format "timestamp" column as dates
df.modify("drop", "useless column") # drops "useless column"
df.interpolate_missing_rows() # interpolate missing rows
df.create_lag(3, "value") # create 3 lag columns from "value"

df.show_all=True # print every row/column of data frame
print(df)

training = TrainingWin(df) # create training window
training.split(["hourOfDay", "dayOfWeek", "monthOfYear", "lag1", "lag2", "lag3"], "value", 0.2)
# split 80/20%

# Tune the hyper-parameters
param_grid = {
    'bootstrap': [True, False],
    'max_depth': [110, None],
    'max_features': [1.0, "sqrt", "log2"],
    'min_samples_leaf': [1, 2],
    'min_samples_split': [2, 6, 8],
    'n_estimators': [1600, 2000]
}
training.tune_hps(param_grid)

# You can also set them manually
# training.hyper_param = {'bootstrap': False, 'max_depth': None, 'max_features': 'sqrt',
# 'min_samples_leaf': 1, 'min_samples_split': 6, 'n_estimators': 1600}

# Generate and save model
training.generate_model()
training.save_model(r"path\to\file.sav")

# You can also load a pre-existing .sav file as a model
# training.load_model(r"C:\Users\omaro\Desktop\Data\Riyadh2-210123971-FebToMay-Model.sav")

# Predict "value" column with "timestamp" column and create an anomaly column with 3 * std
training.predict_csv(r"path\to\file.csv", "timestamp", "value", r"path\to\file.csv", 3)
print(training.predict_list([12, 0, 6, 30, 40, 50])) # prints the prediction for this particular
set of features

# Creates graphing window and graphs heatmap using "value" column
graphing_window = GraphWin(df)
graphing_window.create_graph("heatmap", "value", None)
```