

```

1  # Library imports
2  import numpy as np
3  import pandas as pd
4
5  # Reading original dataset into excel file
6  original_df = pd.read_excel('Original-River-Data.xlsx', usecols='A:I', skiprows=1)
7  river_data = original_df.copy()
8
9  # Renaming Headers
10 new_columns = {'Unnamed: 0': 'Date'}
11 new_columns.update({col: f'{col} MDF (Cumeecs)' for col in river_data.columns[1:5]})
12 new_columns.update({col: f'{col} DRT (mm)' for col in river_data.columns[5:]})
13
14 river_data.rename(
15     columns=new_columns,
16     inplace=True
17 )
18
19 # Converting non-numeric columns to numeric columns
20 river_data['Skip Bridge MDF (Cumeecs)'] = pd.to_numeric(river_data['Skip Bridge MDF (Cumeecs)'], errors='coerce')
21 river_data['Skelton MDF (Cumeecs)'] = pd.to_numeric(river_data['Skelton MDF (Cumeecs)'], errors='coerce')
22 river_data['East Cowton DRT (mm)'] = pd.to_numeric(river_data['East Cowton DRT (mm)'], errors='coerce')
23 river_data['Date'] = pd.to_datetime(river_data['Date'], errors='coerce')
24
25 # Dropping rows with at least 1 null value
26 flow_cols = list(river_data.columns[1:5])
27 rain_cols = list(river_data.columns[5:])
28
29 null_values = river_data.isna().any(axis=1)
30 river_data.dropna(how="any", inplace=True)
31
32 # Dropping rows with rainfall outliers
33 rainfall_outliers = river_data[(river_data[rain_cols] > 400).any(1)]
34 river_data.drop(rainfall_outliers.index, inplace=True)
35
36 # Dropping rows with river flow outliers
37 river_flow_outliers = river_data[(river_data[flow_cols] == 0).any(1)]
38 river_data.drop(river_flow_outliers.index, inplace=True)
39
40 # Exporting cleaned dataset to excel file
41 export_data = river_data.copy()
42 export_data["Date"] = export_data["Date"].astype("string")
43 export_data.to_excel('River-Data-Cleaned.xlsx')
44
45 # Lagging data
46 clean_df = pd.read_excel('River-Data-Cleaned.xlsx')
47 clean_df.drop(["Unnamed: 0"], axis=1, inplace=True)
48
49 lagged_df = pd.DataFrame()
50 lagged_df["Date"] = clean_df["Date"]
51 lagged_df[flow_cols[-1]] = clean_df[flow_cols[-1]]
52
53 ## Lagging rainfall and flow columns by 1 to 3 days
54 for i in range(3):
55     for col in flow_cols:
56         col_name = col.replace("(Cumeecs)", f"(t-{i+1})")
57         lagged_df[col_name] = clean_df[col].shift(i+1)
58
59 for i in range(3):
60     for col in rain_cols:
61         col_name = col.replace("(mm)", f"(t-{i+1})")
62         lagged_df[col_name] = clean_df[col].shift(i+1)
63
64 ## Dropping rows with null values
65 lagged_df[lagged_df.isna().any(axis=1)]
66 lagged_df.dropna(how="any", inplace=True)

```

```

66
67 # Moving averages
68 moving_avg_df = pd.DataFrame()
69 moving_avg_df["Date"] = clean_df["Date"]
70 moving_avg_df[flow_cols[-1]] = clean_df[flow_cols[-1]]
71
72 ## Creating moving averages of between 3 and 7 days for each numerical column
73 for i in range(3, 8):
74     for col in flow_cols:
75         col_name = col.replace("(Cumecs)", f"(MA{i})")
76         moving_avg_df[col_name] = clean_df[col].rolling(i).mean()
77
78 for i in range(3, 8):
79     for col in rain_cols:
80         col_name = col.replace("(mm)", f"(MA{i})")
81         moving_avg_df[col_name] = clean_df[col].rolling(i).mean()
82
83 ## Dropping rows with null values
84 moving_avg_df[moving_avg_df.isna().any(axis=1)]
85 moving_avg_df.dropna(how="any", inplace=True)
86
87 # Lagged moving averages
88 lagged_ma_df = pd.DataFrame()
89 lagged_ma_df["Date"] = moving_avg_df["Date"]
90 lagged_ma_df[flow_cols[-1]] = moving_avg_df[flow_cols[-1]]
91
92 ## Lagging moving averages by 1 day
93 ## lagging them by more than 1 day results in much weaker correlations
94 mdf_cols = list(moving_avg_df.columns[2:22])
95 drt_cols = list(moving_avg_df.columns[22:])
96
97 for col in mdf_cols:
98     col_name = col + f" (t-1)"
99     lagged_ma_df[col_name] = moving_avg_df[col].shift(1)
100
101 for col in drt_cols:
102     col_name = col + f" (t-1)"
103     lagged_ma_df[col_name] = moving_avg_df[col].shift(1)
104
105 ## Dropping rows with null values
106 lagged_ma_df[lagged_ma_df.isna().any(axis=1)]
107 lagged_ma_df.dropna(how="any", inplace=True)
108
109 # Weighted moving averages
110 weighted_ma_df = pd.DataFrame()
111 weighted_ma_df["Date"] = clean_df["Date"]
112 weighted_ma_df[flow_cols[-1]] = clean_df[flow_cols[-1]]
113
114 ## Creating weighted moving averages of between 3 and 7 days
115 ## for each numerical column
116 for i in range(3, 8):
117     for col in flow_cols:
118         col_name = col.replace("(Cumecs)", f"(WMA{i})")
119         weighted_ma_df[col_name] = clean_df[col].ewm(span=i).mean()
120
121 for i in range(3, 8):
122     for col in rain_cols:
123         col_name = col.replace("(mm)", f"(WMA{i})")
124         weighted_ma_df[col_name] = clean_df[col].ewm(span=i).mean()
125
126 # Lagged weighted moving averages
127 lagged_wma_df = pd.DataFrame()
128 lagged_wma_df["Date"] = weighted_ma_df["Date"]
129 lagged_wma_df[flow_cols[-1]] = weighted_ma_df[flow_cols[-1]]
130
131 ## Lagging weighted moving averages
132

```

```

133 w_mdf_cols = list(weighted_ma_df.columns[2:22])
134 w_drt_cols = list(weighted_ma_df.columns[22:])
135
136 for col in w_mdf_cols:
137     col_name = col + f" (t-1)"
138     lagged_wma_df[col_name] = weighted_ma_df[col].shift(1)
139
140 for col in w_drt_cols:
141     col_name = col + f" (t-1)"
142     lagged_wma_df[col_name] = weighted_ma_df[col].shift(1)
143
144 ## Dropping rows with null values
145 lagged_wma_df[lagged_wma_df.isna().any(1)]
146 lagged_wma_df.dropna(how="any", inplace=True)
147
148 # Exporting newly created datasets with lags and moving averages
149 lagged_df.to_excel('River-Data-Lagged.xlsx')
150 moving_avg_df.to_excel('River-Data-MA.xlsx')
151 lagged_ma_df.to_excel('River-Data-MA-Lagged.xlsx')
152 weighted_ma_df.to_excel('River-Data-WMA.xlsx')
153 lagged_wma_df.to_excel('River-Data-WMA-Lagged.xlsx')
154
155 # Utility Functions
156 ## Functions for standardising and unstandardising values
157 def standardise_columns(df, cols):
158     """
159     This function works with dataframes to standardise values
160     in multiple columns to the range [0.1, 0.9]
161     """
162     subset_df = df[cols]
163     subset_df = 0.8 * ((subset_df - subset_df.min()) / (subset_df.max() - subset_df.min())) + 0.1
164     return subset_df
165
166 def unstandardise_columns(df, cols, max_val, min_val):
167     """
168     This function works with numpy arrays to destandardise values
169     in multiple columns
170     """
171     subset_df = df[cols]
172     subset_df = ((subset_df - subset_df.min()) / 0.8) * (max_val - min_val) + min_val
173     return subset_df
174
175 def standardise_value(x, max_val, min_val):
176     """
177     This function works with numpy arrays to standardise values
178     in multiple arrays to the range [0.1, 0.9]
179     """
180     return 0.8 * ((x - min_val) / (max_val - min_val)) + 0.1
181
182 def unstandardise_value(x, max_val, min_val):
183     """
184     This function works with numpy arrays to destandardise values
185     in multiple arrays
186     """
187     return ((x - 0.1) / 0.8) * (max_val - min_val) + min_val
188
189 # Function for building custom feature and target sets from larger datasets
190 def build_feature_set(*datasets):
191     assert len(datasets) > 0, "No data sets entered"
192     datasets = list(datasets)
193     min_rows = min(d.shape[0] for d in datasets)
194
195     for i, ds in enumerate(datasets):
196         datasets[i] = ds.truncate(before=ds.shape[0]-min_rows).reset_index()
197         datasets[i].drop(["index"], axis=1, inplace=True)
198
199

```

```

200 merged_df = datasets[0].iloc[:, :2]
201 for ds in datasets:
202     merged_df = pd.concat([merged_df, ds.iloc[:, 2:]], axis=1)
203
204 merged_cols = list(merged_df.columns)
205 selected_cols = []
206
207 for i in range(0, len(merged_cols), 2):
208     format_str = f"{i+1}) {merged_cols[i]}"
209     if i != len(merged_cols) - 1:
210         second_part = f"{i+2}) {merged_cols[i+1]}"
211         num_spaces = 50 - len(format_str)
212         format_str += num_spaces*" " + second_part
213     print(format_str)
214
215 selected_indices = input("\nSelect columns: ")
216 for index in selected_indices.split(","):
217     if "-" in index:
218         first_i, second_i = index.split("-")
219         selected_cols += merged_cols[int(first_i) - 1: int(second_i)]
220     else:
221         selected_cols.append(merged_cols[int(index) - 1])
222
223 return merged_df[selected_cols]

```