

MACHINE LEARNING

(PREDICTION ALGORITHM)

TOPIC CHOSEN: Time Series Analysis

DATASET LINK: <https://www.kaggle.com/code/andreshg/timeseries-analysis-a-complete-guide/notebook>

DATASET NAME: Aquier_Petrignano.csv

SOURCE: Kaggle

DATASET DETAILS:

In this project I want to deep dive into TimeSeries analysis to show how to review the data, how to preview it and how to engineering.

I also want to explore some of the typical TimeSeries topics such as:

ACF/PACF

ARIMA

Auto-ARIMA

Prophet

Augmented Dickey-Fuller (ADF)

DATA VISUALIZATION:

FEATURES:

Rainfall indicates the quantity of rain falling (mm)

Temperature indicates the temperature (°C)

TARGET:

- **Depth to Groundwater** indicates the groundwater level (m from the ground floor)

```
import pandas as pd
import numpy as np
import datetime as dt
```

```
df = pd.read_csv("Aquifer_Petrignano.csv")
df.head()
```

	Date	Rainfall_Bastia_Umbra	Depth_to_Groundwater_P24	Depth_to_Groundwater_P25	Temperature_Bastia_Umbra	Temperat
0	14/03/2006	NaN	-22.48	-22.18	NaN	
1	15/03/2006	NaN	-22.38	-22.14	NaN	
2	16/03/2006	NaN	-22.25	-22.04	NaN	
3	17/03/2006	NaN	-22.38	-22.04	NaN	
4	18/03/2006	NaN	-22.60	-22.04	NaN	

```
# Remove old rows
df = df[df.Rainfall_Bastia_Umbra.notna()].reset_index(drop=True)
# Remove not usefull columns
df = df.drop(['Depth_to_Groundwater_P24', 'Temperature_Petrignano'], axis=1)
```

```
# Simplify column names
df.columns = ['date', 'rainfall', 'depth_to_groundwater', 'temperature', 'drainage_volume',
'river_hydrometry']
```

```
targets = ['depth_to_groundwater']
features = [feature for feature in df.columns if feature not in targets]
df.head()
```

	date	rainfall	depth_to_groundwater	temperature	drainage_volume	river_hydrometry
0	01/01/2009	0.0	-31.14	5.2	-24530.688	2.4
1	02/01/2009	0.0	-31.11	2.3	-28785.888	2.5
2	03/01/2009	0.0	-31.07	4.4	-25766.208	2.4
3	04/01/2009	0.0	-31.05	0.8	-27919.296	2.4
4	05/01/2009	0.0	-31.01	-1.9	-29854.656	2.3

```
from datetime import datetime, date
```

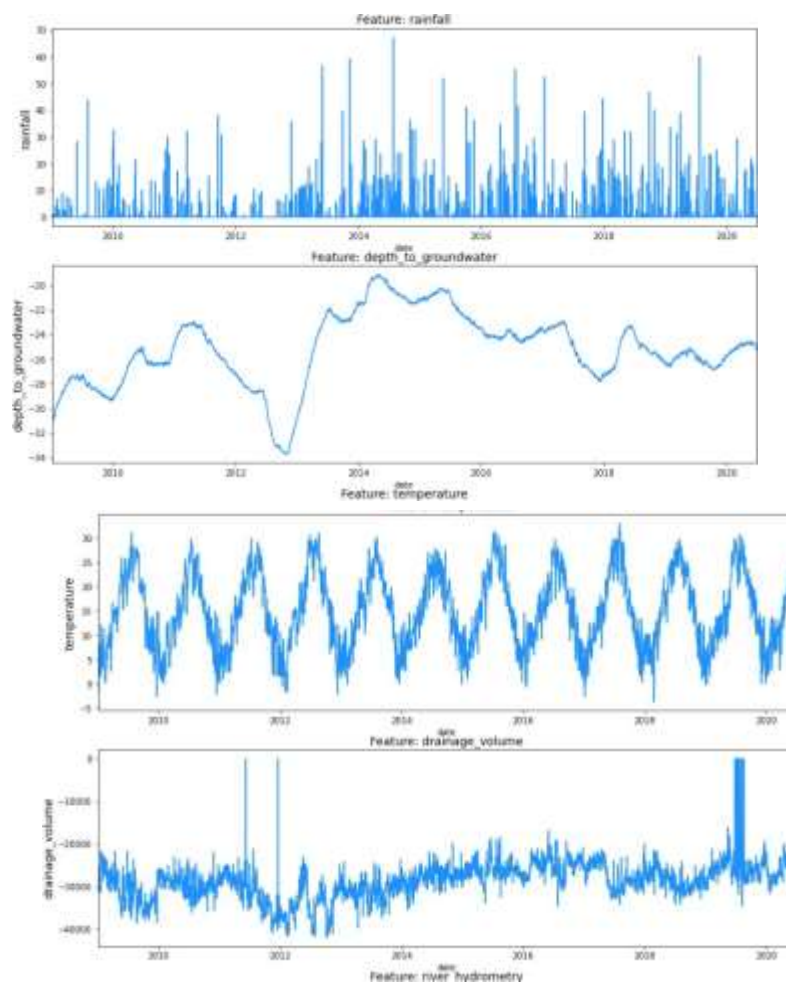
```
df['date'] = pd.to_datetime(df['date'], format = '%d/%m/%Y')
df.head().style.set_properties(subset=['date'], **{'background-color': 'dodgerblue'})
```

	date	rainfall	depth_to_groundwater	temperature	drainage_volume	river_hydrometry
0	2009-01-01 00:00:00	0.000000	-31.140000	5.200000	-24530.688000	2.400000
1	2009-01-02 00:00:00	0.000000	-31.110000	2.300000	-28785.888000	2.500000
2	2009-01-03 00:00:00	0.000000	-31.070000	4.400000	-25766.208000	2.400000
3	2009-01-04 00:00:00	0.000000	-31.050000	0.800000	-27919.296000	2.400000
4	2009-01-05 00:00:00	0.000000	-31.010000	-1.900000	-29854.656000	2.300000

```
import matplotlib.pyplot as plt
import seaborn as sns
f,ax = plt.subplots(nrows=5, ncols=1, figsize=(15, 25))

for i, column in enumerate(df.drop('date', axis=1).columns):
    sns.lineplot(x=df['date'], y=df[column].fillna(method='ffill'), ax=ax[i], color='dodgerblue')
    ax[i].set_title('Feature: {}'.format(column), fontsize=14)
    ax[i].set_ylabel(ylabel=column, fontsize=14)

ax[i].set_xlim([date(2009, 1, 1), date(2020, 6, 30)])
```



DATA PREPROCESSING:

```
df = df.sort_values(by='date')
```

```
# Check time intervals
```

```
df['delta'] = df['date'] - df['date'].shift(1)
```

```
df[['date', 'delta']].head()
```

	date	delta
0	2009-01-01	NaT
1	2009-01-02	1 days
2	2009-01-03	1 days
3	2009-01-04	1 days
4	2009-01-05	1 days

```
df['delta'].sum(), df['delta'].count()
```

```
(Timedelta('4198 days 00:00:00'), 4198)
```

HANDLE MISSING DATA:

```
df = df.drop('delta', axis=1)
```

```
df.isna().sum()
```

```
date          0
rainfall      0
depth_to_groundwater  27
temperature   0
drainage_volume  1
river_hydrometry  0
dtype: int64
```

```
f, ax = plt.subplots(nrows=2, ncols=1, figsize=(15, 15))
```

```
old_hydrometry = df['river_hydrometry'].copy()
```

```
df['river_hydrometry'] = df['river_hydrometry'].replace(0, np.nan)
```

```
sns.lineplot(x=df['date'], y=old_hydrometry, ax=ax[0], color='darkorange', label='original')
```

```
sns.lineplot(x=df['date'], y=df['river_hydrometry'].fillna(np.inf), ax=ax[0], color='dodgerblue',  
label='modified')
```

```
ax[0].set_title('Feature: Hydrometry', fontsize=14)
```

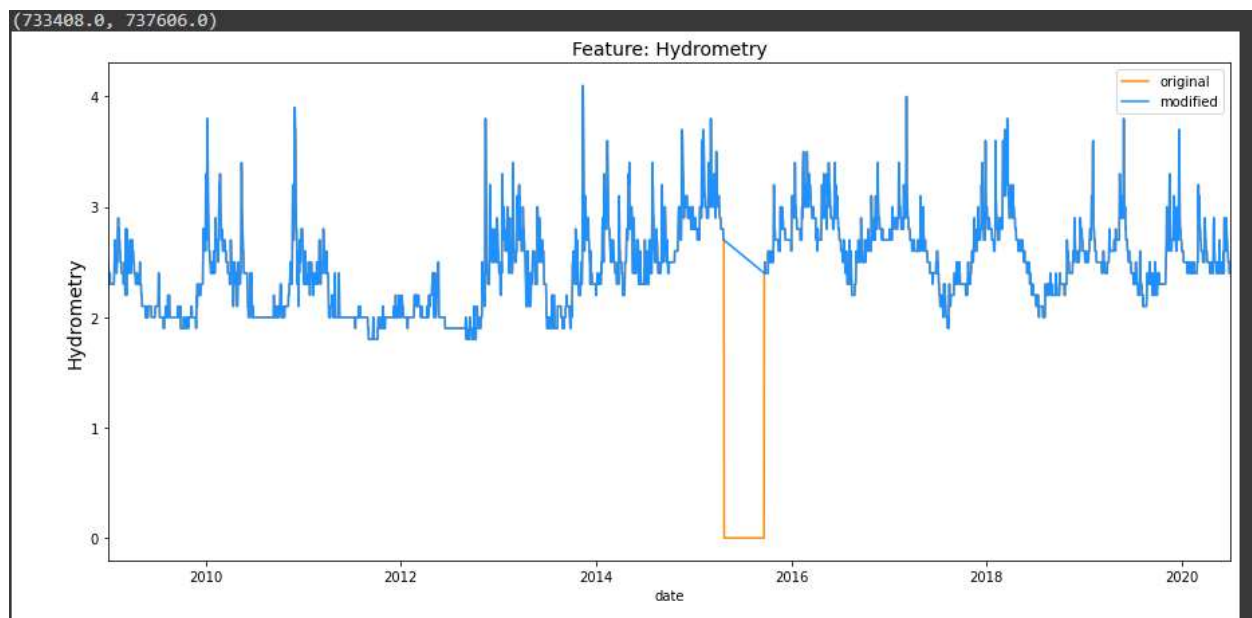
```

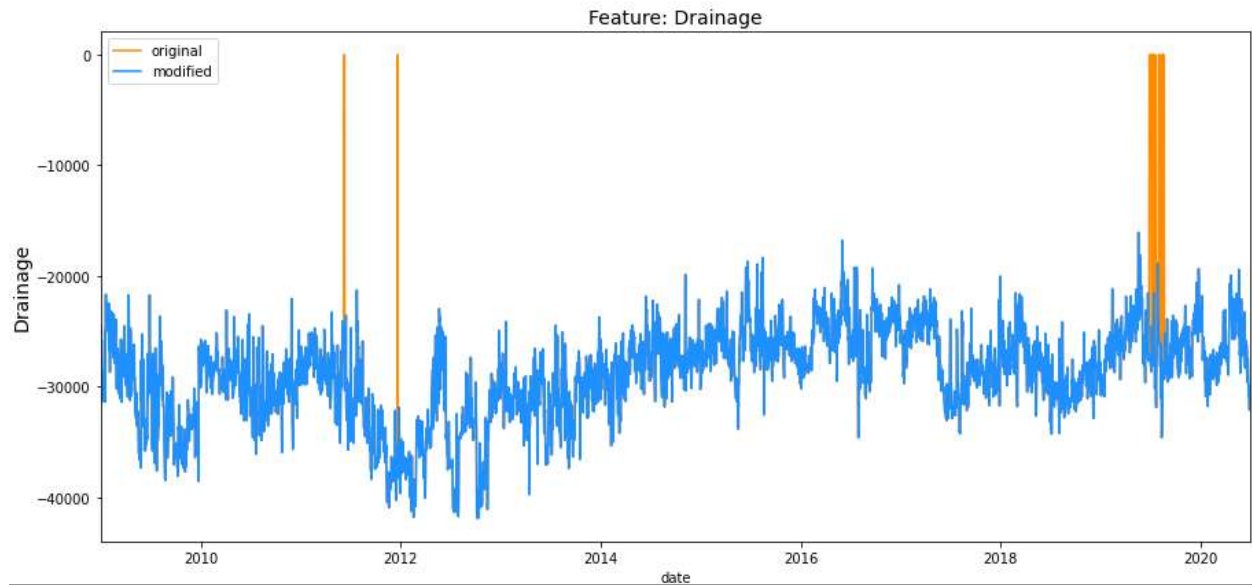
ax[0].set_ylabel(ylabel='Hydrometry', fontsize=14)
ax[0].set_xlim([date(2009, 1, 1), date(2020, 6, 30)])

old_drainage = df['drainage_volume'].copy()
df['drainage_volume'] = df['drainage_volume'].replace(0, np.nan)

sns.lineplot(x=df['date'], y=old_drainage, ax=ax[1], color='darkorange', label='original')
sns.lineplot(x=df['date'], y=df['drainage_volume'].fillna(np.inf), ax=ax[1], color='dodgerblue',
label='modified')
ax[1].set_title('Feature: Drainage', fontsize=14)
ax[1].set_ylabel(ylabel='Drainage', fontsize=14)
ax[1].set_xlim([date(2009, 1, 1), date(2020, 6, 30)])

```

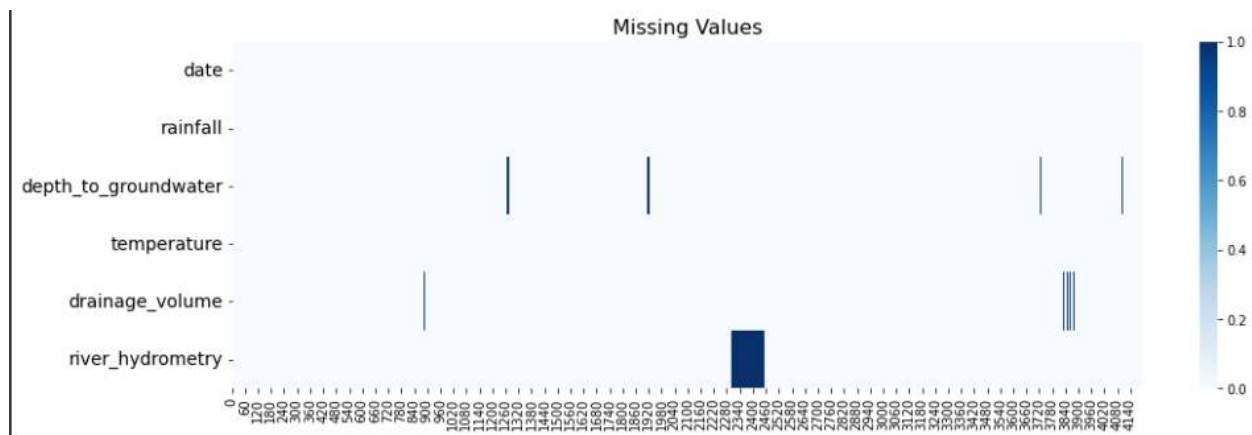




```
f, ax = plt.subplots(nrows=1, ncols=1, figsize=(16,5))
```

```
sns.heatmap(df.T.isna(), cmap='Blues')
ax.set_title('Missing Values', fontsize=16)
```

```
for tick in ax.yaxis.get_major_ticks():
    tick.label.set_fontsize(14)
plt.show()
```



```

f, ax = plt.subplots(nrows=4, ncols=1, figsize=(15, 12))

sns.lineplot(x=df['date'], y=df['drainage_volume'].fillna(0), ax=ax[0], color='darkorange', label =
'modified')
sns.lineplot(x=df['date'], y=df['drainage_volume'].fillna(np.inf), ax=ax[0], color='dodgerblue', label =
'original')
ax[0].set_title('Fill NaN with 0', fontsize=14)
ax[0].set_ylabel(ylabel='Volume C10 Petignano', fontsize=14)

mean_drainage = df['drainage_volume'].mean()
sns.lineplot(x=df['date'], y=df['drainage_volume'].fillna(mean_drainage), ax=ax[1], color='darkorange',
label = 'modified')
sns.lineplot(x=df['date'], y=df['drainage_volume'].fillna(np.inf), ax=ax[1], color='dodgerblue', label =
'original')
ax[1].set_title(f'Fill NaN with Mean Value ({mean_drainage:.0f})', fontsize=14)
ax[1].set_ylabel(ylabel='Volume C10 Petignano', fontsize=14)

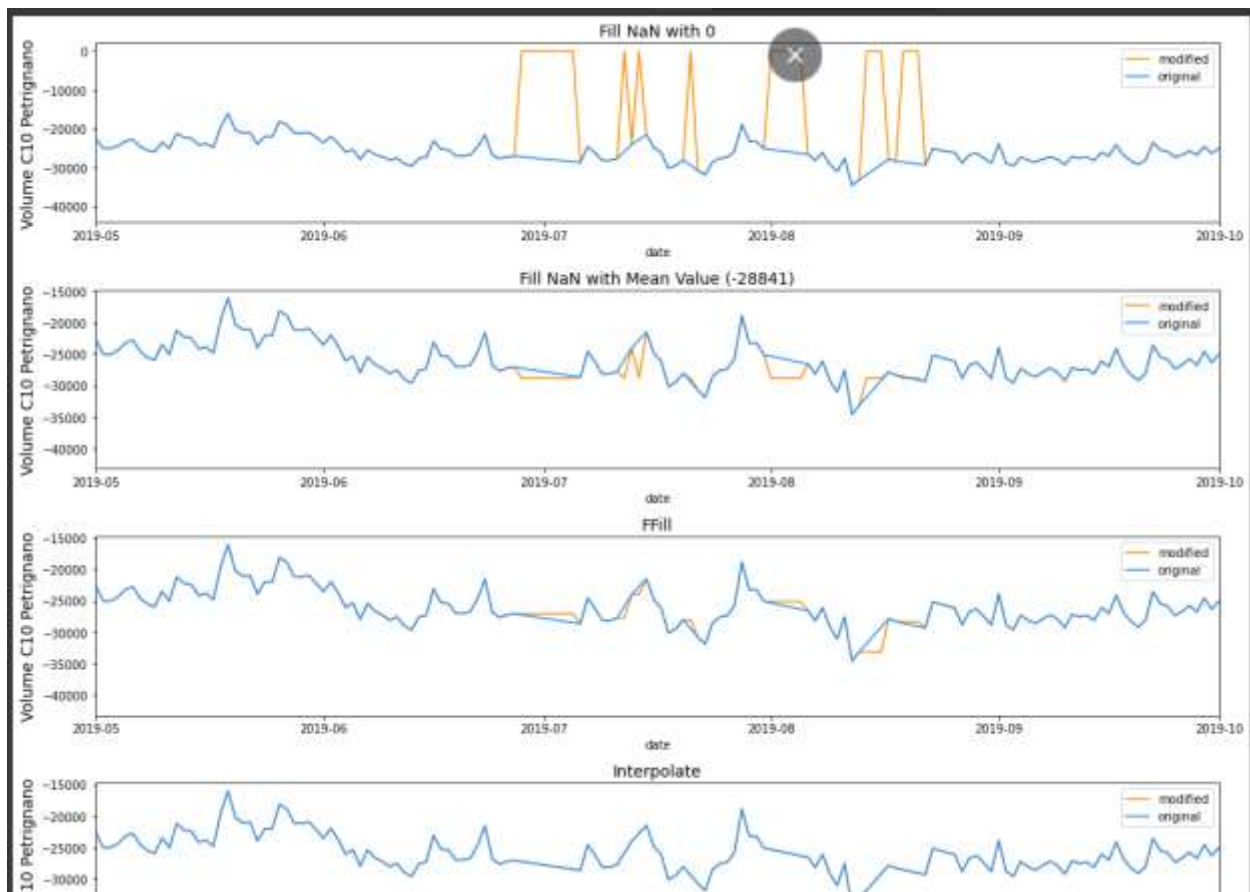
sns.lineplot(x=df['date'], y=df['drainage_volume'].ffill(), ax=ax[2], color='darkorange', label = 'modified')
sns.lineplot(x=df['date'], y=df['drainage_volume'].fillna(np.inf), ax=ax[2], color='dodgerblue', label =
'original')
ax[2].set_title(f'FFill', fontsize=14)
ax[2].set_ylabel(ylabel='Volume C10 Petignano', fontsize=14)

sns.lineplot(x=df['date'], y=df['drainage_volume'].interpolate(), ax=ax[3], color='darkorange', label =
'modified')
sns.lineplot(x=df['date'], y=df['drainage_volume'].fillna(np.inf), ax=ax[3], color='dodgerblue', label =
'original')
ax[3].set_title(f'Interpolate', fontsize=14)
ax[3].set_ylabel(ylabel='Volume C10 Petignano', fontsize=14)

for i in range(4):
    ax[i].set_xlim([date(2019, 5, 1), date(2019, 10, 1)])

plt.tight_layout()
plt.show()

```



```
df['drainage_volume'] = df['drainage_volume'].interpolate()
df['river_hydrometry'] = df['river_hydrometry'].interpolate()
df['depth_to_groundwater'] = df['depth_to_groundwater'].interpolate()
```

Smoothing data:

Resampling can provide additional information on the data. There are two types of resampling:

- **Upsampling** is when the frequency of samples is increased (e.g. days to hours)
- **Downsampling** is when the frequency of samples is decreased (e.g. days to weeks)

```
fig, ax = plt.subplots(ncols=2, nrows=3, sharex=True, figsize=(16,12))
```

```
sns.lineplot(df['date'], df['drainage_volume'], color='dodgerblue', ax=ax[0, 0])
ax[0, 0].set_title('Drainage Volume', fontsize=14)
```

```
resampled_df = df[['date', 'drainage_volume']].resample('7D', on='date').sum().reset_index(drop=False)
sns.lineplot(resampled_df['date'], resampled_df['drainage_volume'], color='dodgerblue', ax=ax[1, 0])
ax[1, 0].set_title('Weekly Drainage Volume', fontsize=14)
```

```
resampled_df = df[['date', 'drainage_volume']].resample('M', on='date').sum().reset_index(drop=False)
```



```
sns.lineplot(resampled_df['date'], resampled_df['drainage_volume'], color='dodgerblue', ax=ax[2, 0])
ax[2, 0].set_title('Monthly Drainage Volume', fontsize=14)
```

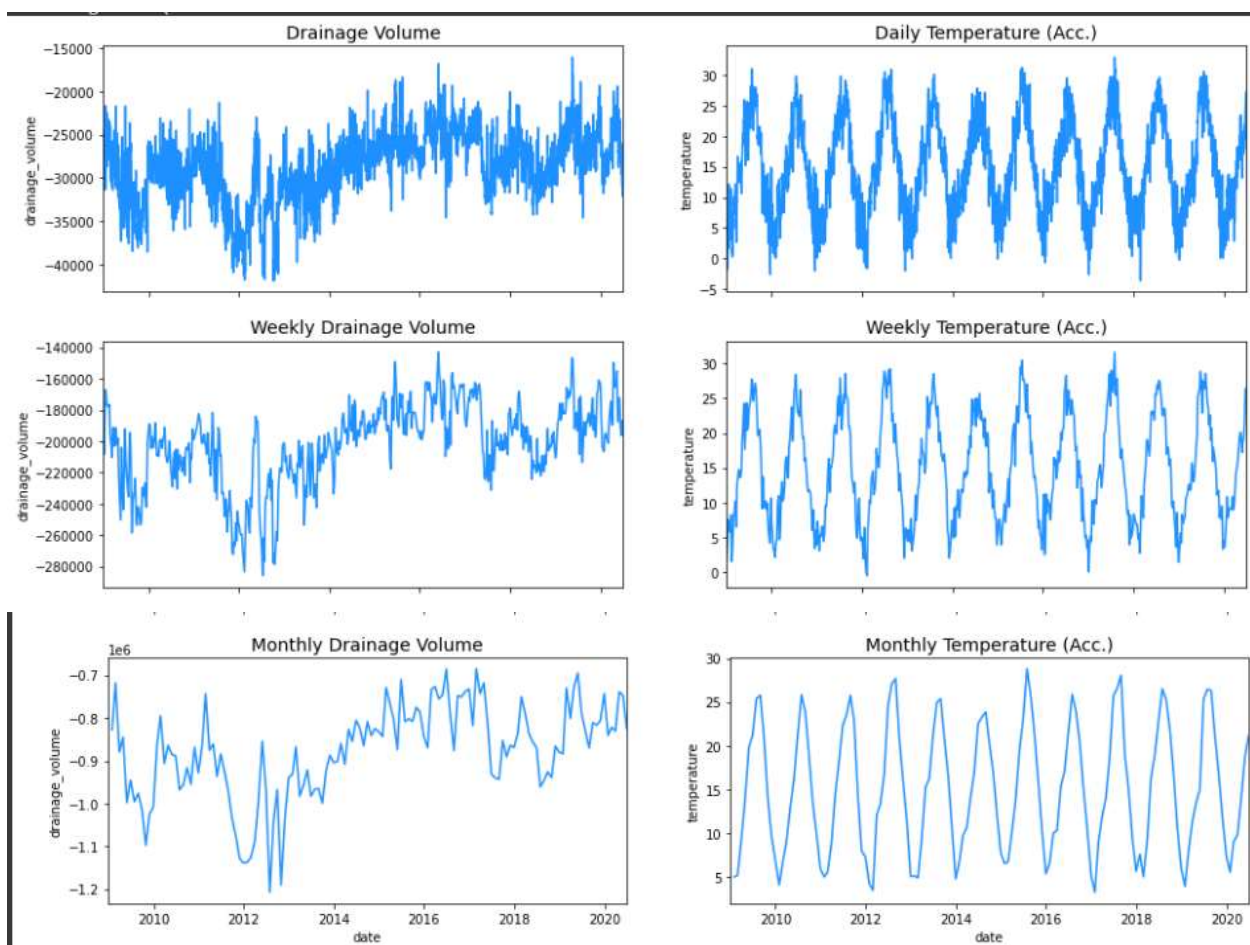
```
for i in range(3):
    ax[i, 0].set_xlim([date(2009, 1, 1), date(2020, 6, 30)])
```

```
sns.lineplot(df['date'], df['temperature'], color='dodgerblue', ax=ax[0, 1])
ax[0, 1].set_title('Daily Temperature (Acc.)', fontsize=14)
```

```
resampled_df = df[['date', 'temperature']].resample('7D', on='date').mean().reset_index(drop=False)
sns.lineplot(resampled_df['date'], resampled_df['temperature'], color='dodgerblue', ax=ax[1, 1])
ax[1, 1].set_title('Weekly Temperature (Acc.)', fontsize=14)
```

```
resampled_df = df[['date', 'temperature']].resample('M', on='date').mean().reset_index(drop=False)
sns.lineplot(resampled_df['date'], resampled_df['temperature'], color='dodgerblue', ax=ax[2, 1])
ax[2, 1].set_title('Monthly Temperature (Acc.)', fontsize=14)
```

```
for i in range(3):
    ax[i, 1].set_xlim([date(2009, 1, 1), date(2020, 6, 30)])
plt.show()
```



```

downsample = df[['date',
                 'depth_to_groundwater',
                 'temperature',
                 'drainage_volume',
                 'river_hydrometry',
                 'rainfall'
                ]].resample('7D', on='date').mean().reset_index(drop=False)

```

```
df = downsample.copy()
```

STATIONARITY:

Some time-series models, such as such as ARIMA, assume that the underlying data is stationary. Stationarity describes that the time-series has

- constant mean and mean is not time-dependent
- constant variance and variance is not time-dependent
- constant covariance and covariance is not time-dependent

```
rolling_window = 52
```

```
f, ax = plt.subplots(nrows=2, ncols=1, figsize=(15, 12))
```

```

sns.lineplot(x=df['date'], y=df['drainage_volume'], ax=ax[0], color='dodgerblue')
sns.lineplot(x=df['date'], y=df['drainage_volume'].rolling(rolling_window).mean(), ax=ax[0],
color='black', label='rolling mean')
sns.lineplot(x=df['date'], y=df['drainage_volume'].rolling(rolling_window).std(), ax=ax[0], color='orange',
label='rolling std')
ax[0].set_title('Depth to Groundwater: Non-stationary \nnon-constant mean & non-constant variance',
fontsize=14)
ax[0].set_ylabel(ylabel='Drainage Volume', fontsize=14)
ax[0].set_xlim([date(2009, 1, 1), date(2020, 6, 30)])

```

```

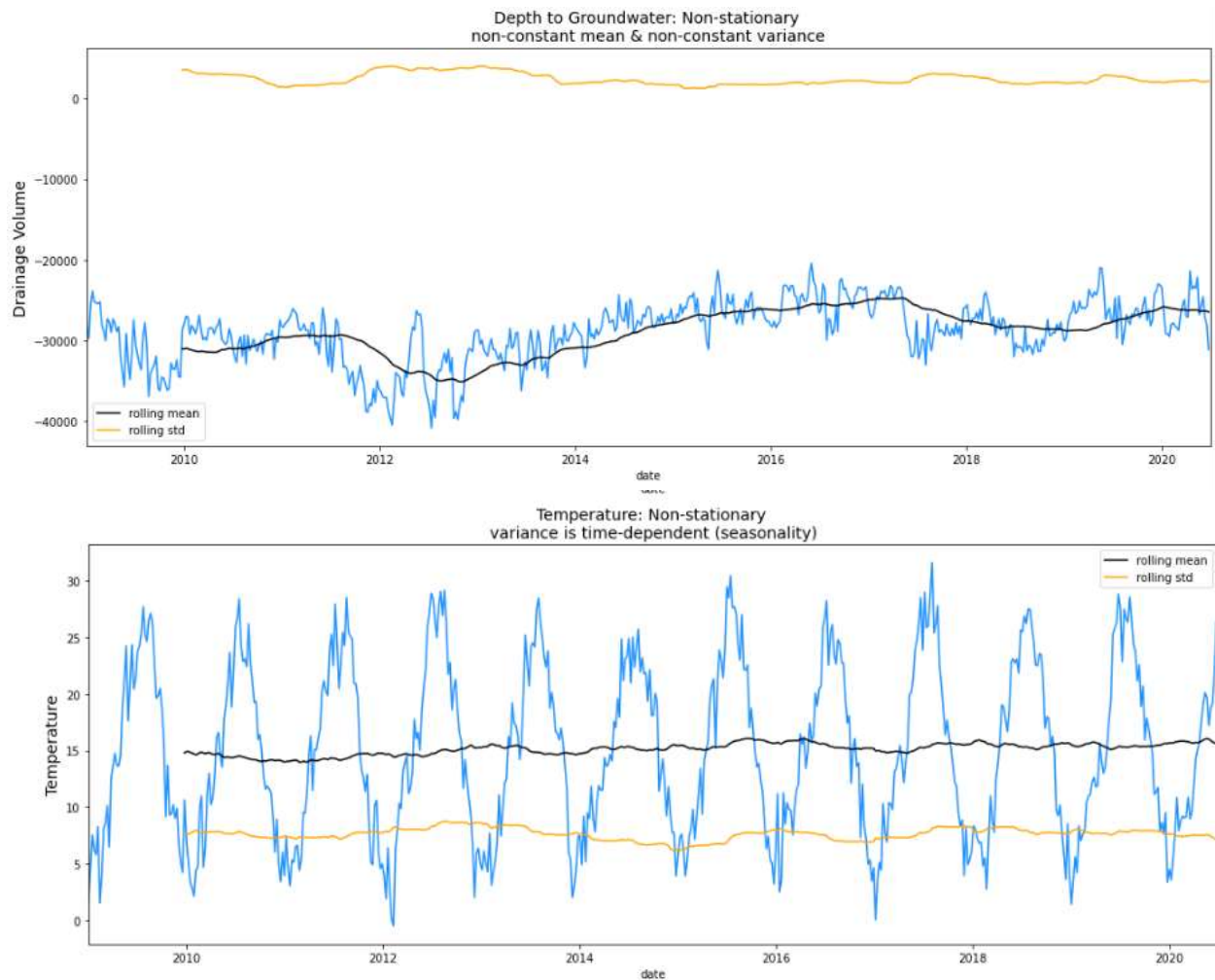
sns.lineplot(x=df['date'], y=df['temperature'], ax=ax[1], color='dodgerblue')
sns.lineplot(x=df['date'], y=df['temperature'].rolling(rolling_window).mean(), ax=ax[1], color='black',
label='rolling mean')
sns.lineplot(x=df['date'], y=df['temperature'].rolling(rolling_window).std(), ax=ax[1], color='orange',
label='rolling std')
ax[1].set_title('Temperature: Non-stationary \nvariance is time-dependent (seasonality)', fontsize=14)
ax[1].set_ylabel(ylabel='Temperature', fontsize=14)
ax[1].set_xlim([date(2009, 1, 1), date(2020, 6, 30)])

```

```

plt.tight_layout()
plt.show()

```



Feature engineering:

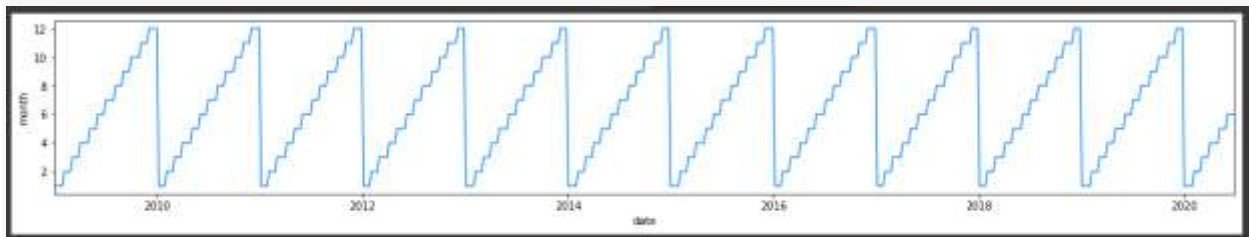
```
df['year'] = pd.DatetimeIndex(df['date']).year
df['month'] = pd.DatetimeIndex(df['date']).month
df['day'] = pd.DatetimeIndex(df['date']).day
df['day_of_year'] = pd.DatetimeIndex(df['date']).dayofyear
df['week_of_year'] = pd.DatetimeIndex(df['date']).weekofyear
df['quarter'] = pd.DatetimeIndex(df['date']).quarter
df['season'] = df['month'] % 12 // 3 + 1

df[['date', 'year', 'month', 'day', 'day_of_year', 'week_of_year', 'quarter', 'season']].head()
```

	date	year	month	day	day_of_year	week_of_year	quarter	season
0	2009-01-01	2009	1	1	1	1	1	1
1	2009-01-08	2009	1	8	8	2	1	1
2	2009-01-15	2009	1	15	15	3	1	1
3	2009-01-22	2009	1	22	22	4	1	1
4	2009-01-29	2009	1	29	29	5	1	1

```
f, ax = plt.subplots(nrows=1, ncols=1, figsize=(20, 3))

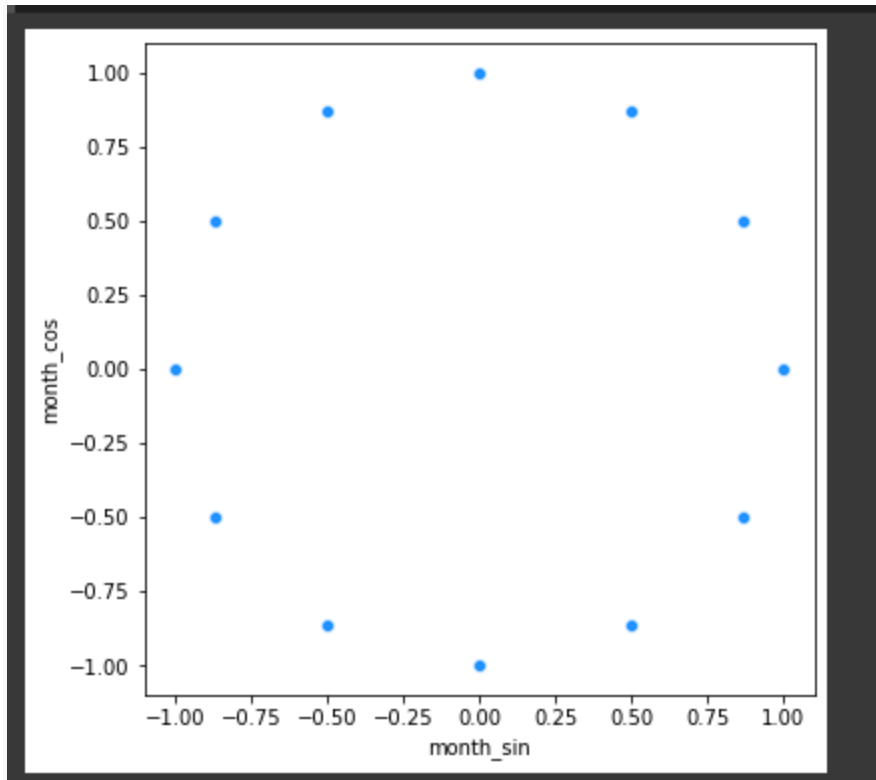
sns.lineplot(x=df['date'], y=df['month'], color='dodgerblue')
ax.set_xlim([date(2009, 1, 1), date(2020, 6, 30)])
plt.show()
```



```
month_in_year = 12
df['month_sin'] = np.sin(2*np.pi*df['month']/month_in_year)
df['month_cos'] = np.cos(2*np.pi*df['month']/month_in_year)

f, ax = plt.subplots(nrows=1, ncols=1, figsize=(6, 6))

sns.scatterplot(x=df.month_sin, y=df.month_cos, color='dodgerblue')
plt.show()
```



TIME SERIES DECOMPOSITION:

Time series decomposition involves thinking of a series as a combination of level, trend, seasonality, and noise components.

These components are defined as follows:

- **Level:** The average value in the series.
- **Trend:** The increasing or decreasing value in the series.
- **Seasonality:** The repeating short-term cycle in the series.
- **Noise:** The random variation in the series.

It is helpful to think of the components as combining either additively or multiplicatively:

- **Additive:** $y(t) = \text{Level} + \text{Trend} + \text{Seasonality} + \text{Noise}$
- **Multiplicative:** $y(t) = \text{Level} * \text{Trend} * \text{Seasonality} * \text{Noise}$

```
from statsmodels.tsa.seasonal import seasonal_decompose
```

```

core_columns = [
    'rainfall', 'temperature', 'drainage_volume',
    'river_hydrometry', 'depth_to_groundwater'
]

for column in core_columns:
    decomp = seasonal_decompose(df[column], period=52, model='additive', extrapolate_trend='freq')
    df[f"{column}_trend"] = decomp.trend
    df[f"{column}_seasonal"] = decomp.seasonal

```

```

fig, ax = plt.subplots(ncols=2, nrows=4, sharex=True, figsize=(16,8))

```

```

for i, column in enumerate(['temperature', 'depth_to_groundwater']):

```

```

    res = seasonal_decompose(df[column], freq=52, model='additive', extrapolate_trend='freq')

```

```

    ax[0,i].set_title('Decomposition of {}'.format(column), fontsize=16)
    res.observed.plot(ax=ax[0,i], legend=False, color='dodgerblue')
    ax[0,i].set_ylabel('Observed', fontsize=14)

```

```

    res.trend.plot(ax=ax[1,i], legend=False, color='dodgerblue')
    ax[1,i].set_ylabel('Trend', fontsize=14)

```

```

    res.seasonal.plot(ax=ax[2,i], legend=False, color='dodgerblue')
    ax[2,i].set_ylabel('Seasonal', fontsize=14)

```

```

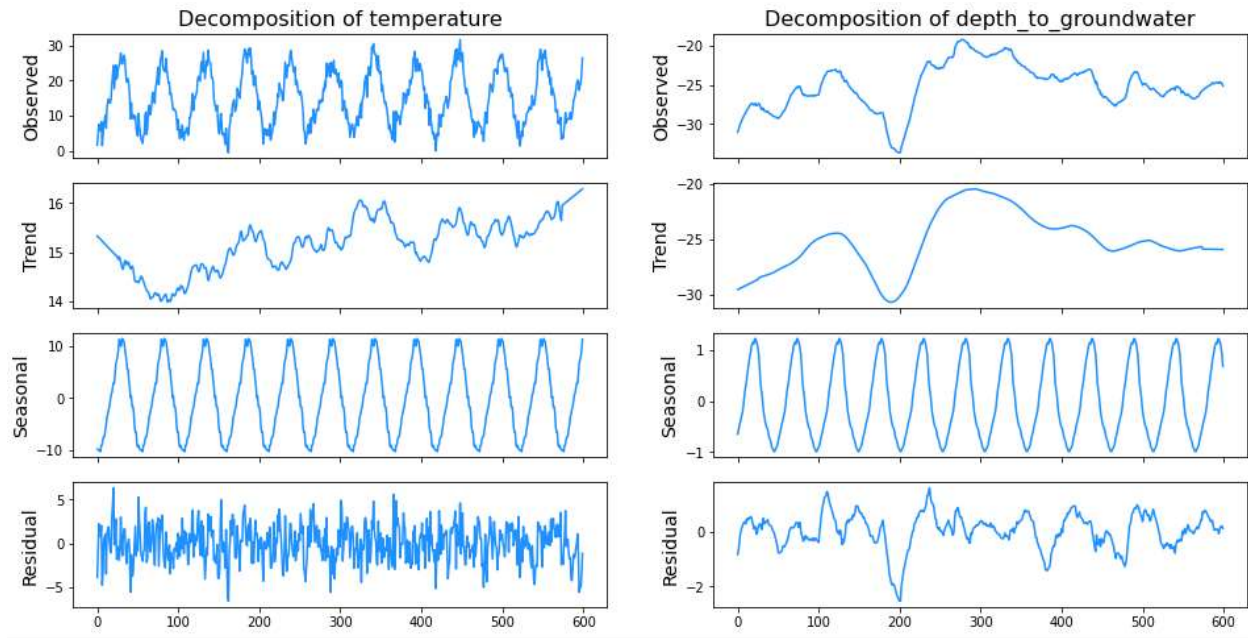
    res.resid.plot(ax=ax[3,i], legend=False, color='dodgerblue')
    ax[3,i].set_ylabel('Residual', fontsize=14)

```

```

plt.show()

```



LAG :

We want to calculate each variable with a `shift()` (lag) to compare the correlation with the other variables.

```
weeks_in_month = 4

for column in core_columns:
    df[f'{column}_seasonal_shift_b_2m'] = df[f'{column}_seasonal'].shift(-2 * weeks_in_month)
    df[f'{column}_seasonal_shift_b_1m'] = df[f'{column}_seasonal'].shift(-1 * weeks_in_month)
    df[f'{column}_seasonal_shift_1m'] = df[f'{column}_seasonal'].shift(1 * weeks_in_month)
    df[f'{column}_seasonal_shift_2m'] = df[f'{column}_seasonal'].shift(2 * weeks_in_month)
    df[f'{column}_seasonal_shift_3m'] = df[f'{column}_seasonal'].shift(3 * weeks_in_month)
```

EXPLORATORY DATA ANALYSIS:

Now, we are going to plot the data and try to extract some knowledge.

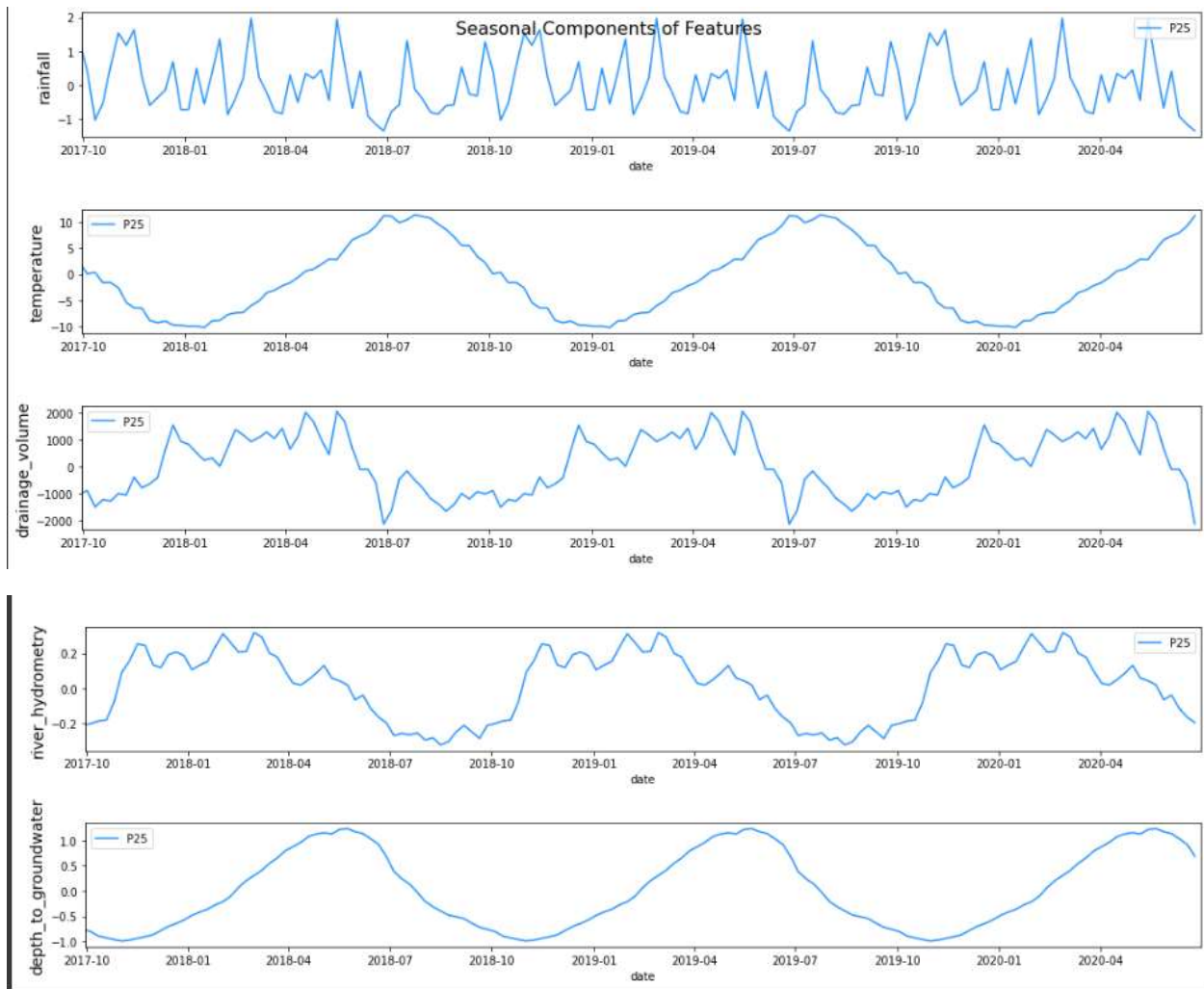
```
f, ax = plt.subplots(nrows=5, ncols=1, figsize=(15, 12))
f.suptitle('Seasonal Components of Features', fontsize=16)
```

```

for i, column in enumerate(core_columns):
    sns.lineplot(x=df['date'], y=df[column + '_seasonal'], ax=ax[i], color='dodge
rblue', label='P25')
    ax[i].set_ylabel(ylabel=column, fontsize=14)
    ax[i].set_xlim([date(2017, 9, 30), date(2020, 6, 30)])

plt.tight_layout()
plt.show()

```



```

f, ax = plt.subplots(nrows=1, ncols=2, figsize=(16, 8))

corrmat = df[core_columns].corr()

sns.heatmap(corrmat, annot=True, vmin=-1, vmax=1, cmap='coolwarm_r', ax=ax[0])
ax[0].set_title('Correlation Matrix of Core Features', fontsize=16)

shifted_cols = [
    'depth_to_groundwater_seasonal',

```



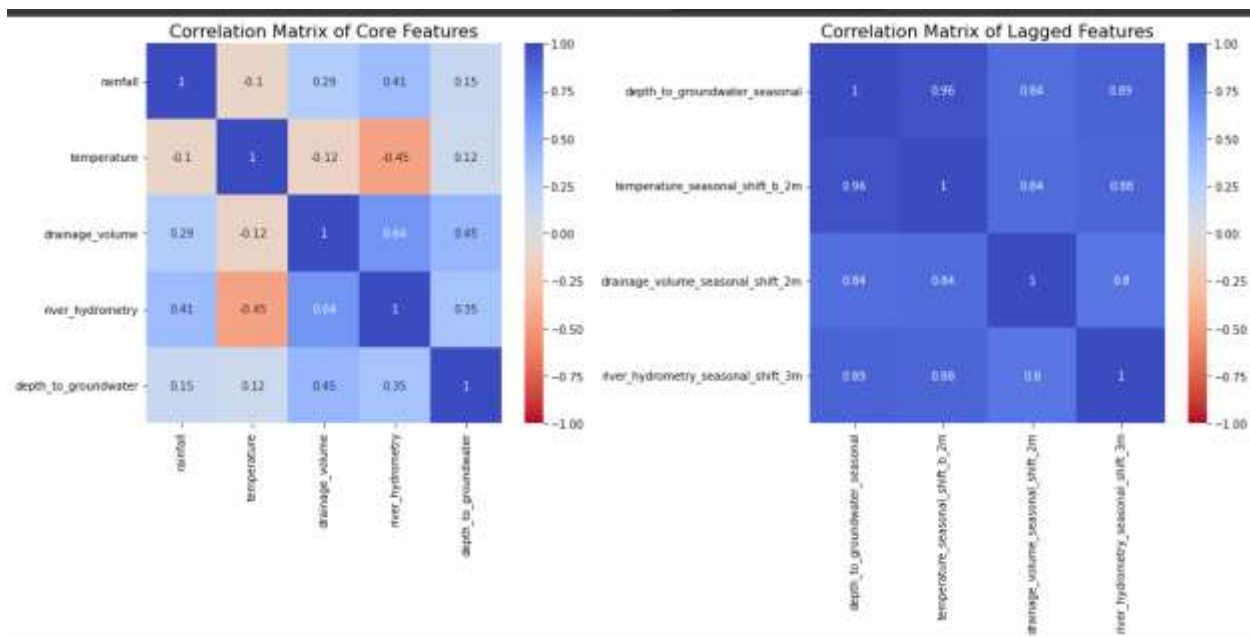
```

        'temperature_seasonal_shift_b_2m',
        'drainage_volume_seasonal_shift_2m',
        'river_hydrometry_seasonal_shift_3m'
    ]
    corrmatrix = df[shifted_cols].corr()

    sns.heatmap(corrmatrix, annot=True, vmin=-1, vmax=1, cmap='coolwarm_r', ax=ax[1])
    ax[1].set_title('Correlation Matrix of Lagged Features', fontsize=16)

plt.tight_layout()
plt.show()

```



MODELING:

Time series can be either univariate or multivariate:

- **Univariate** time series only has a single time-dependent variable.
- **Multivariate** time series have a multiple time-dependent variable.

But, first of all we are going to see how does cross-validation technic works in TimeSeries Analysis

```

from sklearn.model_selection import TimeSeriesSplit

N_SPLITS = 3

X = df['date']

```

```

y = df['depth_to_groundwater']
folds = TimeSeriesSplit(n_splits=N_SPLITS)

```

```

f, ax = plt.subplots(nrows=N_SPLITS, ncols=2, figsize=(16, 9))

for i, (train_index, valid_index) in enumerate(folds.split(X)):
    X_train, X_valid = X[train_index], X[valid_index]
    y_train, y_valid = y[train_index], y[valid_index]

    sns.lineplot(
        x=X_train,
        y=y_train,
        ax=ax[i,0],
        color='dodgerblue',
        label='train'
    )
    sns.lineplot(
        x=X_train[len(X_train) - len(X_valid):(len(X_train) - len(X_valid) + len(X_val
lid))],
        y=y_train[len(X_train) - len(X_valid):(len(X_train) - len(X_valid) + len(X_val
lid))],
        ax=ax[i,1],
        color='dodgerblue',
        label='train'
    )

    for j in range(2):
        sns.lineplot(x= X_valid, y= y_valid, ax=ax[i, j], color='darkorange', label='
validation')
        ax[i, 0].set_title(f"Rolling Window with Adjusting Training Size (Split {i+1})",
        fontsize=16)
        ax[i, 1].set_title(f"Rolling Window with Constant Training Size (Split {i+1})", f
ontsize=16)

for i in range(N_SPLITS):
    ax[i, 0].set_xlim([date(2009, 1, 1), date(2020, 6, 30)])
    ax[i, 1].set_xlim([date(2009, 1, 1), date(2020, 6, 30)])

plt.tight_layout()
plt.show()

```

