


## ✓ CIND820 Capstone Project

# Integrating Predictive Analytics and Anomaly Detection for Cyanobacterial Bloom Forecasting

- ✓ *This notebook details the development of predictive analytics for cyanobacterial blooms using ARIMA, LSTM, and anomaly detection models. It includes data preparation, exploratory analysis, and initial modeling insights.*

```
1 # Mount Google Drive to permit colab Notebook access to project datasets
2
3 from google.colab import drive
4 drive.mount('/content/drive')
```

 Mounted at /content/drive

## ✓ 1. Data Collection and Preprocessing

### ✓ 1.1 Import Necessary Libraries

Import libraries that will be required for data analysis, visualization, and handling.

```
1 # Import necessary libraries for data manipulation and visualization.
2
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7
```

### ✓ 1.2 Source Identification and Data Acquisition

Define the file paths to the datasets that will be used in the analysis and in assessing the project research questions.

```
1 # Specify file paths for the datasets
2
3 bloom_indices_path = '/content/drive/MyDrive/CIND820/LakeErie_Daily_BloomIndices.csv'
4 water_quality_path = '/content/drive/MyDrive/CIND820/lake_erie_habs_field_sampling_results_2012_2018.csv'
5
```

### ✓ 1.3 Load Datasets

Read the Bloom Indices and Water Quality datasets into pandas DataFrames for manipulation.

```
1 # Load the Bloom Indices dataset
2 bloom_data = pd.read_csv(bloom_indices_path, encoding='ISO-8859-1')
3
4 # Load the Water Quality dataset
5 water_quality_data = pd.read_csv(water_quality_path, encoding='ISO-8859-1')
6
```

### ✓ 1.4 Data Inspection

Get a preliminary understanding of the data structure and content.

```
1 # Inspect the first few rows of the Bloom Indices dataset
2 print("Bloom Data Sample:")
3 print(bloom_data.head())
4
5 # Print a line of underscores to demarcate the outputs
6 print("_ " * 70)
7
8 # Inspect the first few rows of the Water Quality dataset
9 print("\nWater Quality Data Sample:")
10 print(water_quality_data.head())
11
```



```

3          41.8339          -83.3640          NaN
4          41.7625          -83.3286          NaN

Wave Height (ft) ... Soluble Reactive Phosphorus (µg P/L) \
0          NaN ... 3.23
1          NaN ... 2.02
2          NaN ... 2.97
3          NaN ... 3.02
4          NaN ... 5.15

Ammonia (µg N/L) Nitrate + Nitrite (mg N/L) Urea (µg N/L) \
0          10.22          0.46          NaN
1          28.79          0.32          NaN
2          35.7          0.52          NaN
3          11.28          0.51          NaN
4          14.68          0.47          NaN

Particulate Organic Carbon (mg/L) Particulate Organic Nitrogen (mg/L) \
0          NaN          NaN
1          0.34          0.05
2          0.40          0.08
3          0.49          0.12
4          0.50          0.07

Dissolved Organic Carbon (mg/L) \
0          2.88
1          1.80
2          4.35
3          2.53
4          2.80

Colored Dissolved Organic Material absorbance (m-1) at 400nm \
0          NaN
1          NaN
2          NaN
3          NaN
4          NaN

Total Suspended Solids (mg/L) Volatile Suspended Solids (mg/L)
0          3.47          1.11
1          3.15          1.01
2          2.90          0.98
3          4.32          1.06
4          21.42          2.40

```

```
[5 rows x 36 columns]
```

## ✓ 1.5 Data Cleaning and Preparation

### ✓ 1.5.1 Parsing Dates

Convert date columns to datetime objects for time-series analysis.

2. Convert bloom\_data Dates to Match water\_quality\_data Use the pd.to\_datetime function to parse bloom\_data dates into yyyy-mm-dd:

```

1 # Inspect the current date formats and data types
2 print("Initial Date Format and Data Type:")
3 print(bloom_data['Date'].head())
4 print(water_quality_data['Date'].head())
5
6 print("\nData Types:")
7 print(f"Bloom Data Date Type: {bloom_data['Date'].dtype}")
8 print(f"Water Quality Data Date Type: {water_quality_data['Date'].dtype}")
9
10 # Convert Date columns to datetime
11 bloom_data['Date'] = pd.to_datetime(bloom_data['Date'], format='%Y-%m-%d', errors='coerce')
12 water_quality_data['Date'] = pd.to_datetime(water_quality_data['Date'], format='%m/%d/%Y', errors='coerce')
13
14 # Check for invalid dates (after conversion)
15 print("\nPost Conversion - Check for Null Dates:")
16 print(bloom_data['Date'].isnull().sum())
17 print(water_quality_data['Date'].isnull().sum())
18
19 # Verify the new data types
20 print("\nPost Conversion - Data Types:")
21 print(f"Bloom Data Date Type: {bloom_data['Date'].dtype}")
22 print(f"Water Quality Data Date Type: {water_quality_data['Date'].dtype}")
23
24 # Preview the converted dates
25 print("\nConverted Date Formats:")
26 print(bloom_data['Date'].head())
27 print(water_quality_data['Date'].head())
28

```

Initial Date Format and Data Type:

```

0    2002-06-02
1    2002-06-15
2    2002-06-18
3    2002-06-23
4    2002-06-24
Name: Date, dtype: object
0    5/15/2012
1    5/15/2012
2    5/15/2012
3    5/15/2012
4    5/31/2012
Name: Date, dtype: object

```

Data Types:

Bloom Data Date Type: object

```

Water Quality Data Date Type: object

Post Conversion - Check for Null Dates:
0
0

Post Conversion - Data Types:
Bloom Data Date Type: datetime64[ns]
Water Quality Data Date Type: datetime64[ns]

Converted Date Formats:
0 2002-06-02
1 2002-06-15
2 2002-06-18
3 2002-06-23
4 2002-06-24
Name: Date, dtype: datetime64[ns]
0 2012-05-15
1 2012-05-15
2 2012-05-15
3 2012-05-15
4 2012-05-31
Name: Date, dtype: datetime64[ns]

```

## ✓ 1.5.2 Handling Column Names

Ensure consistency in column names for easier data manipulation. Standardize Column Names Normalize column names by stripping whitespace and other hidden characters.

2. Standardize Column Names Normalize column names by stripping whitespace and other hidden characters:

```

1 bloom_data.columns = bloom_data.columns.str.strip().str.replace('\n', '').str.replace('\r', '').str.replace(' ', ' ')
2 water_quality_data.columns = water_quality_data.columns.str.strip().str.replace('\n', '').str.replace('\r', '').str.replace(' ', ' ')
3
4 print("Normalized column names in bloom_data:", bloom_data.columns.tolist())
5 print("Normalized column names in water_quality_data:", water_quality_data.columns.tolist())
6

```

```

➤ Normalized column names in bloom_data: ['Date', 'Lake', 'Satellite Sensor', 'Bloom Extent (KM2)', 'Bloom Extent (% of Lake Area)', 'Bloom Intensity (µg/L)',
Normalized column names in water_quality_data: ['Date', 'Site', 'Station Depth (m)', 'Sample Depth (m)', 'Sample Depth (category)', 'Local Time (Eastern Tim

```

## ✓ 1.5.3 Handling Missing Values

Address missing data to prevent errors in analysis.

```

1 # Check for missing values in Bloom Data
2 print("\nMissing values in Bloom Data:")
3 print(bloom_data.isnull().sum())
4

```



```

Missing values in Bloom Data:
Date                                0
Lake                                0
Satellite Sensor                    0
Bloom Extent (KM2)                  0
Bloom Extent (% of Lake Area)       0
Bloom Intensity (µg/L)              0
Bloom Severity (µg/L km2)          0
Valid Pixels (% of Lake Area)      0
dtype: int64

```

```

1 # Check for missing values in Water Quality Data
2 print("\nMissing values in Water Quality Data:")
3 print(water_quality_data.isnull().sum())
4

```



```

Missing values in Water Quality Data:
Date                                0
Site                                0
Station Depth (m)                   371
Sample Depth (m)                     17
Sample Depth (category)              0
Local Time (Eastern Time Zone)      363
Latitude (decimal deg)              376
Longitude (decimal deg)             376
Wind speed (knots)                  650
Wave Height (ft)                    655
Sky                                  366
Secchi Depth (m)                     405
Sample Temperature (°C)             1188
CTD Temperature (°C)                 136
CTD Specific Conductivity (µS/cm)    133
CTD Beam Attenuation (m-1)          151
CTD Transmission (%)                 151
CTD Dissolved Oxygen (mg/L)         132
CTD Photosynthetically Active Radiation (µE/m2/s) 134
Turbidity (NTU)                     424
Particulate Microcystin (µg/L)       19
Dissolved Microcystin (µg/L)         251
Extracted Phycocyanin (µg/L)         4
Chlorophyll_a                        5
Total Phosphorus (µg P/L)            371
Total Dissolved Phosphorus (µg P/L)  369
Soluble Reactive Phosphorus (µg P/L) 367
Ammonia (µg N/L)                    366

```

```

Nitrate + Nitrite (mg N/L)      367
Urea (µg N/L)                  925
Particulate Organic Carbon (mg/L) 371
Particulate Organic Nitrogen (mg/L) 371
Dissolved Organic Carbon (mg/L)  524
Colored Dissolved Organic Material absorbance (m-1) at 400nm 568
Total Suspended Solids (mg/L)    383
Volatile Suspended Solids (mg/L) 381
dtype: int64

```

```

1 # Handle missing values in Water Quality Data
2 # For numeric columns, we can fill missing values with mean or median if appropriate
3 numeric_columns = water_quality_data.select_dtypes(include=[np.number]).columns.tolist()
4 water_quality_data[numeric_columns] = water_quality_data[numeric_columns].fillna(method='ffill')
5

```

## ✓ 1.6 Data Integration

### ✓ 1.6.1 Aligning Datasets Temporally

Combine datasets based on the date to facilitate joint analysis.

Final Check Rerun the merge operation after ensuring the Date columns are correctly formatted and aligned:

```

1 # Merge datasets on 'Date'
2 merged_data = pd.merge(bloom_data, water_quality_data, on='Date', how='inner')
3

```

```

1 merged_data = pd.merge(bloom_data, water_quality_data, on='Date', how='inner')
2 print(merged_data.head())
3

```

```

[ ] Bloom Extent (% of Lake Area)  Bloom Intensity (µg/L)  \
0                                0.51                49.26
1                                0.51                49.26
2                                0.51                49.26
3                                0.51                49.26
4                                0.51                49.26

```

```

4          0703.03          99.94 WE8

Station Depth (m) ... Soluble Reactive Phosphorus (µg P/L) \
0          5.50 ...          1.4
1          NaN ...          NaN
2          8.71 ...          0.9
3          3.10 ...          2.7
4          4.67 ...          3.4

Ammonia (µg N/L) Nitrate + Nitrite (mg N/L) Urea (µg N/L) \
0          17.9          0.47          8.96
1          NaN          NaN          NaN
2          56.1          0.34          10.85
3          7.4          1.73          8.37
4          38.9          1.02          16.21

Particulate Organic Carbon (mg/L) Particulate Organic Nitrogen (mg/L) \
0          0.40          0.05
1          NaN          NaN
2          0.88          0.16
3          2.37          0.37
4          0.44          0.07

Dissolved Organic Carbon (mg/L) \
0          2.30
1          NaN
2          1.94
3          4.53
4          2.83

Colored Dissolved Organic Material absorbance (m-1) at 400nm \
0          0.37
1          NaN
2          0.22
3          1.68
4          0.68

Total Suspended Solids (mg/L) Volatile Suspended Solids (mg/L)
0          2.02          0.20
1          NaN          NaN
2          12.27          1.12
3          24.53          3.52
4          2.42          0.50

```

```
[5 rows x 43 columns]
```

```

1 # Having already created the `merged_data`...
2
3 # Convert the merged data to a DataFrame
4 merged_df = pd.DataFrame(merged_data)
5
6 # Display the entire DataFrame in a scrollable table using Pandas
7 from IPython.display import display
8

```



```

9 # Display the DataFrame
10 display(merged_df)
11

```



	Date	Lake	Satellite Sensor	Bloom Extent (KM2)	Bloom Extent (% of Lake Area)	Bloom Intensity (µg/L)	Bloom Severity (µg/L km2)	Valid Pixels (% of Lake Area)	Site	Station Depth (m)	...	Soluble Reactive Phosphorus (µg P/L)	Ammonia (µg N/L)	Nitrate + Nitrite (mg N/L)	Urea (µg N/L)	Particulate Organic Carbon (mg/L)	Particulate Organic Nitrogen (mg/L)	...
0	2016-06-13	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE2	5.50	...	1.4	17.9	0.47	8.96	0.40	0.05	...
1	2016-06-13	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE2	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	...
2	2016-06-13	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE4	8.71	...	0.9	56.1	0.34	10.85	0.88	0.16	...
3	2016-06-13	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE6	3.10	...	2.7	7.4	1.73	8.37	2.37	0.37	...
4	2016-06-13	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE8	4.67	...	3.4	38.9	1.02	16.21	0.44	0.07	...
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
610	2018-10-01	Lake Erie	OLCI-S3	164.07	0.61	25.82	4236.66	99.91	WE13	9.07	...	6.17	1.72	0.27	NaN	0.50	0.10	...
611	2018-10-01	Lake Erie	OLCI-S3	164.07	0.61	25.82	4236.66	99.91	WE13	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	...
612	2018-10-01	Lake Erie	OLCI-S3	164.07	0.61	25.82	4236.66	99.91	WE16	6.30	...	18.33	22.69	0.11	NaN	0.72	0.12	...
613	2018-10-09	Lake Erie	OLCI-S3	157.86	0.59	24.85	3923.31	99.90	WE6	NaN	...	44.14	78.41	0.63	NaN	1.19	0.21	...
614	2018-10-09	Lake Erie	OLCI-S3	157.86	0.59	24.85	3923.31	99.90	WE9	NaN	...	59.19	145.34	0.98	NaN	1.45	0.27	...

615 rows × 43 columns



```

1 from google.colab import data_table
2 data_table.DataTable(merged_df)
3
4

```

Warning: Total number of columns (43) exceeds max\_columns (20) limiting to first (20) columns.

1 to 25 of 615 entries

Filter



index	Date	Lake	Satellite Sensor	Bloom Extent (KM2)	Bloom Extent (% of Lake Area)	Bloom Intensity (µg/L)	Bloom Severity (µg/L km2)	Valid Pixels (% of Lake Area)	Site	Stat
0	2016-06-13 00:00:00	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE2	
1	2016-06-13 00:00:00	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE2	
2	2016-06-13 00:00:00	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE4	
3	2016-06-13 00:00:00	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE6	
4	2016-06-13 00:00:00	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE8	
5	2016-06-13 00:00:00	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE8	
6	2016-06-13 00:00:00	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE9	
7	2016-06-13 00:00:00	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE12	
8	2016-06-13 00:00:00	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE12	
9	2016-06-13 00:00:00	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE13	
10	2016-06-13 00:00:00	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE13	
11	2016-06-13 00:00:00	Lake Erie	OLCI-S3	136.08	0.51	49.26	6703.63	99.94	WE15	
12	2016-06-27 00:00:00	Lake Erie	OLCI-S3	154.53	0.58	48.97	7567.06	99.93	WE2	
13	2016-06-27 00:00:00	Lake Erie	OLCI-S3	154.53	0.58	48.97	7567.06	99.93	WE2	
14	2016-06-27 00:00:00	Lake Erie	OLCI-S3	154.53	0.58	48.97	7567.06	99.93	WE4	
15	2016-06-27 00:00:00	Lake Erie	OLCI-S3	154.53	0.58	48.97	7567.06	99.93	WE6	

16	2016-06-27 00:00:00	Lake Erie	OLCI-S3	154.53	0.58	48.97	7567.06	99.93	WE8
17	2016-06-27 00:00:00	Lake Erie	OLCI-S3	154.53	0.58	48.97	7567.06	99.93	WE8
18	2016-06-27 00:00:00	Lake Erie	OLCI-S3	154.53	0.58	48.97	7567.06	99.93	WE9
19	2016-06-27 00:00:00	Lake Erie	OLCI-S3	154.53	0.58	48.97	7567.06	99.93	WE12
20	2016-06-27 00:00:00	Lake Erie	OLCI-S3	154.53	0.58	48.97	7567.06	99.93	WE12
21	2016-06-27 00:00:00	Lake Erie	OLCI-S3	154.53	0.58	48.97	7567.06	99.93	WE13
22	2016-06-27 00:00:00	Lake Erie	OLCI-S3	154.53	0.58	48.97	7567.06	99.93	WE13
23	2016-06-27 00:00:00	Lake Erie	OLCI-S3	154.53	0.58	48.97	7567.06	99.93	WE15
24	2016-07-05 00:00:00	Lake Erie	OLCI-S3	149.49	0.56	51.05	7631.85	99.94	WE2

◀		▶
Show	25 ▼ per page	
		1 2 10 20 25
◀		▶

```

1 # Define the file path for saving the merged dataframe
2 output_file_path = "/content/drive/MyDrive/CIND820/LakeErie_Merged_Data.csv"
3
4 # Save the merged dataframe to the specified location as a CSV file
5 merged_df.to_csv(output_file_path, index=False)
6
7 # Confirm the file has been saved
8 print(f"Merged DataFrame saved to: {output_file_path}")
9

```

➔ Merged DataFrame saved to: /content/drive/MyDrive/CIND820/LakeErie\_Merged\_Data.csv

### Summary of Section 1

- Handled missing values and ensured consistency across datasets.
- Ensured compatibility for downstream alignment and analysis steps.

## ✓ 2. Exploratory Data Analysis (EDA)

### ✓ 2.1 Statistical Summary and Visualization

#### ✓ 2.1.1 Descriptive Statistics

Provide a statistical overview of the datasets.

```

1 # Summary statistics for Bloom Data
2 print("\nBloom Data Statistics:")
3 print(bloom_data.describe())
4
5 # Summary statistics for Water Quality Data
6 print("\nWater Quality Data Statistics:")
7 print(water_quality_data.describe())
8

```



Bloom Data Statistics:

	Date	Bloom Extent (KM2)	\
count	1369	1369.000000	
mean	2010-10-09 18:57:03.813002240	443.157239	
min	2002-06-02 00:00:00	0.270000	
25%	2006-07-16 00:00:00	174.960000	
50%	2009-08-26 00:00:00	256.320000	

75%	2016-09-10 00:00:00	469.980000
max	2018-10-31 00:00:00	5256.810000
std	NaN	594.930353

	Bloom Extent (% of Lake Area)	Bloom Intensity ( $\mu\text{g/L}$ ) \
count	1369.000000	1369.000000
mean	1.660767	33.203046
min	0.000000	11.190000
25%	0.660000	25.380000
50%	0.960000	31.390000
75%	1.760000	40.320000
max	19.700000	65.620000
std	2.229544	10.133595

	Bloom Severity ( $\mu\text{g/L km}^2$ )	Valid Pixels (% of Lake Area)
count	1369.000000	1369.000000
mean	13417.767480	96.355566
min	3.020000	1.010000
25%	7021.980000	99.290000
50%	8969.150000	99.890000
75%	11618.340000	99.980000
max	176410.750000	100.000000
std	17402.623445	12.535047

#### Water Quality Data Statistics:

	Date	Station Depth (m)	Sample Depth (m) \
count	1244	873.000000	1227.000000
mean	2016-04-14 17:03:16.784566016	5.643860	1.930929
min	2012-05-15 00:00:00	1.900000	0.000000
25%	2015-06-22 00:00:00	4.000000	0.750000
50%	2016-07-18 00:00:00	5.330000	0.750000
75%	2017-08-21 00:00:00	7.800000	3.300000
max	2018-10-09 00:00:00	10.860000	8.700000
std	NaN	2.202444	2.187762

	Latitude (decimal deg)	Longitude (decimal deg)	Wave Height (ft) \
count	868.000000	868.000000	589.000000
mean	41.748882	-83.272040	0.879677
min	41.012700	-83.708400	0.020000
25%	41.705875	-83.363600	0.480000
50%	41.743000	-83.328350	0.750000
75%	41.826300	-83.193375	1.170000
max	42.002000	-82.698300	2.970000
std	0.080505	0.118646	0.539996

	Sample Temperature ( $^{\circ}\text{C}$ )	CTD Temperature ( $^{\circ}\text{C}$ ) \
count	56.000000	1108.000000
mean	23.823214	21.970036
min	16.300000	2.400000
25%	21.700000	20.300000

```

1 # Convert the bloom_data summary statistics output to a DataFrame
2 bloom_data_summary_statistics_df = pd.DataFrame(bloom_data.describe())
3

```

```

4 # Display the DataFrame using Pandas
5 from IPython.display import display
6
7 # Display the DataFrame
8 display(bloom_data_summary_statistics_df)
9
10 # Define the file path for saving the merged dataframe
11 output_file_path = "/content/drive/MyDrive/CIND820/bloom_data_summary_statistics_Output.csv"
12
13 # Save the merged dataframe to the specified location as a CSV file
14 bloom_data_summary_statistics_df.to_csv(output_file_path, index=False)

```



	Bloom Extent (KM2)	Bloom Extent (% of Lake Area)	Bloom Intensity (ug/L)	Bloom Severity (ug/L km2)	Valid Pixels (% of Lake Area)
<b>count</b>	1369.000000	1369.000000	1369.000000	1369.000000	1369.000000
<b>mean</b>	443.157239	1.660767	33.203046	13417.767480	96.355566
<b>std</b>	594.930353	2.229544	10.133595	17402.623445	12.535047
<b>min</b>	0.270000	0.000000	11.190000	3.020000	1.010000
<b>25%</b>	174.960000	0.660000	25.380000	7021.980000	99.290000
<b>50%</b>	256.320000	0.960000	31.390000	8969.150000	99.890000
<b>75%</b>	469.980000	1.760000	40.320000	11618.340000	99.980000
<b>max</b>	5256.810000	10.700000	65.620000	176410.750000	100.000000



Next steps:

[Generate code with bloom\\_data\\_summary\\_statistics\\_df](#)
[View recommended plots](#)
[New interactive sheet](#)

```

1 # Convert the water_quality_data summary statistics output to a DataFrame
2 water_quality_data_summary_statistics_df = pd.DataFrame(water_quality_data.describe())
3
4 # Display the DataFrame using Pandas
5 from IPython.display import display
6
7 # Display the DataFrame
8 display(water_quality_data_summary_statistics_df)
9
10 # Define the file path for saving the merged dataframe
11 output_file_path = "/content/drive/MyDrive/CIND820/water_quality_data_summary_statistics_Output.csv"
12
13 # Save the merged dataframe to the specified location as a CSV file
14 water_quality_data_summary_statistics_df.to_csv(output_file_path, index=False)

```



	Date	Station Depth (m)	Sample Depth (m)	Latitude (decimal deg)	Longitude (decimal deg)	Wave Height (ft)	Sample Temperature (°C)	CTD Temperature (°C)	CTD Specific Conductivity (μS/cm)	CTD Photosynthetically Active Radiation (μE/m2/s)	Turbidity (NTU)
<b>count</b>	1244	873.000000	1227.000000	868.000000	868.000000	589.000000	56.000000	1108.000000	1111.000000	1110.000000	820.000000
<b>mean</b>	2016-04-14 17:03:16.784566016	5.643860	1.930929	41.748882	-83.272040	0.879677	23.823214	21.970036	289.774887	324.566315	19.161780
<b>min</b>	2012-05-15 00:00:00	1.900000	0.000000	41.012700	-83.708400	0.020000	16.300000	2.400000	1.500000	0.000000	0.680000
<b>25%</b>	2015-06-22 00:00:00	4.000000	0.750000	41.705875	-83.363600	0.480000	21.700000	20.300000	249.300000	14.000000	4.627500
<b>50%</b>	2016-07-18 00:00:00	5.330000	0.750000	41.743000	-83.328350	0.750000	24.200000	22.900000	274.000000	147.370000	9.710000
<b>75%</b>	2017-08-21 00:00:00	7.800000	3.300000	41.826300	-83.193375	1.170000	26.050000	24.500000	314.800000	460.372500	19.300000
<b>max</b>	2018-10-09 00:00:00	10.860000	8.700000	42.002000	-82.698300	2.970000	28.300000	29.700000	855.900000	1902.000000	1148.000000
<b>std</b>	NaN	2.202444	2.187762	0.080505	0.118646	0.539996	2.671261	3.635855	64.250139	418.007239	54.735476

Next steps:

[Generate code with water\\_quality\\_data\\_summary\\_statistics\\_df](#)
[View recommended plots](#)
[New interactive sheet](#)

## 2. Exploratory Data Analysis (EDA)


Inspect Column Names to verify the exact column names. Standardize Column Names to remove unnecessary spaces and special characters like "μ" which may not render correctly and to simplify access for downstream analyses.

```
1 print(bloom_data.columns.tolist())
2
```

```
['Date', 'Lake', 'Satellite Sensor', 'Bloom Extent (KM2)', 'Bloom Extent (% of Lake Area)', 'Bloom Intensity (μg/L)', 'Bloom Severity (μg/L km2)', 'Valid Pi
```

```
1 bloom_data.columns = bloom_data.columns.str.strip().str.replace('μ', 'u').str.replace('µ', 'u')
2 print(bloom_data.columns.tolist())
3
```

➤ ['Date', 'Lake', 'Satellite Sensor', 'Bloom Extent (KM2)', 'Bloom Extent (% of Lake Area)', 'Bloom Intensity (ug/L)', 'Bloom Severity (ug/L km2)', 'Valid Pi



## ✓ 2.1.2 Visualization of Distributions

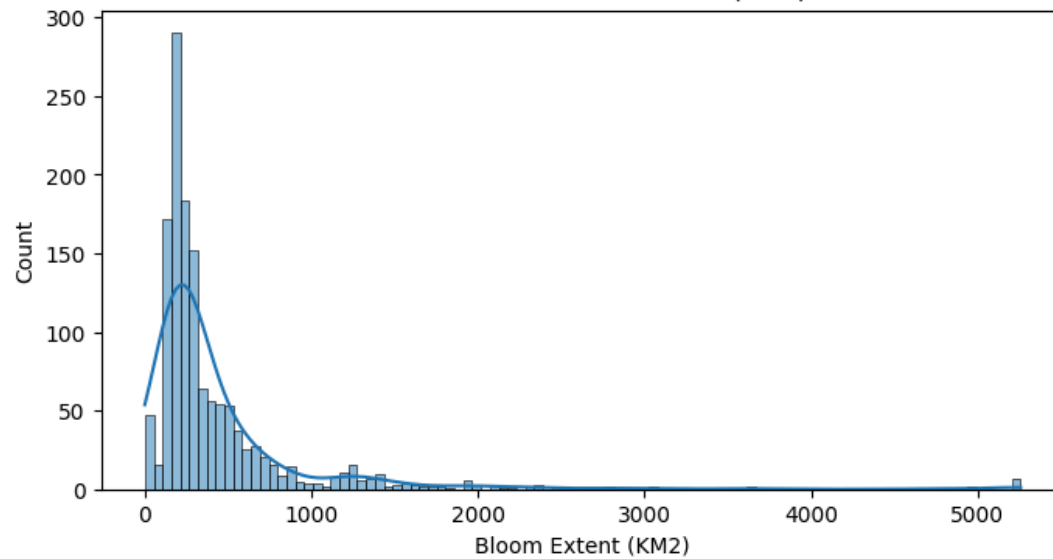
Visualize the distribution of key bloom indices to understand data variability.

```
1 # Plot histograms of key Bloom Indices
2 bloom_features = ['Bloom Extent (KM2)', 'Bloom Intensity (ug/L)', 'Bloom Severity (ug/L km2)']
3
4 for feature in bloom_features:
5     plt.figure(figsize=(8, 4))
6     sns.histplot(bloom_data[feature], kde=True)
7     plt.title(f'Distribution of {feature}')
8     plt.show()
9
```

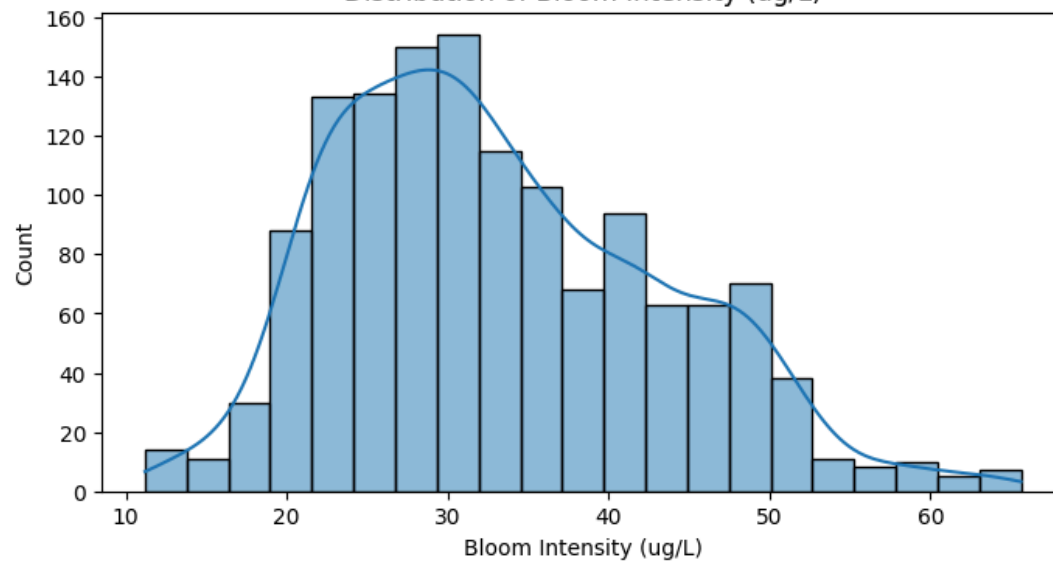




Distribution of Bloom Extent (KM2)

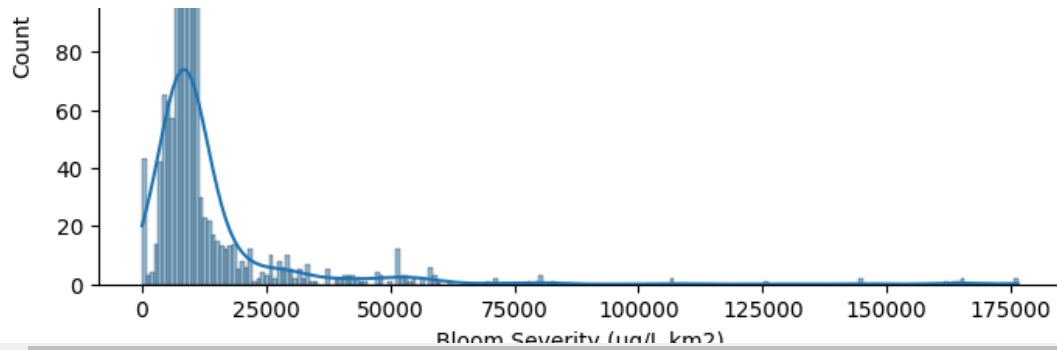


Distribution of Bloom Intensity (ug/L)



Distribution of Bloom Severity (ug/L km2)





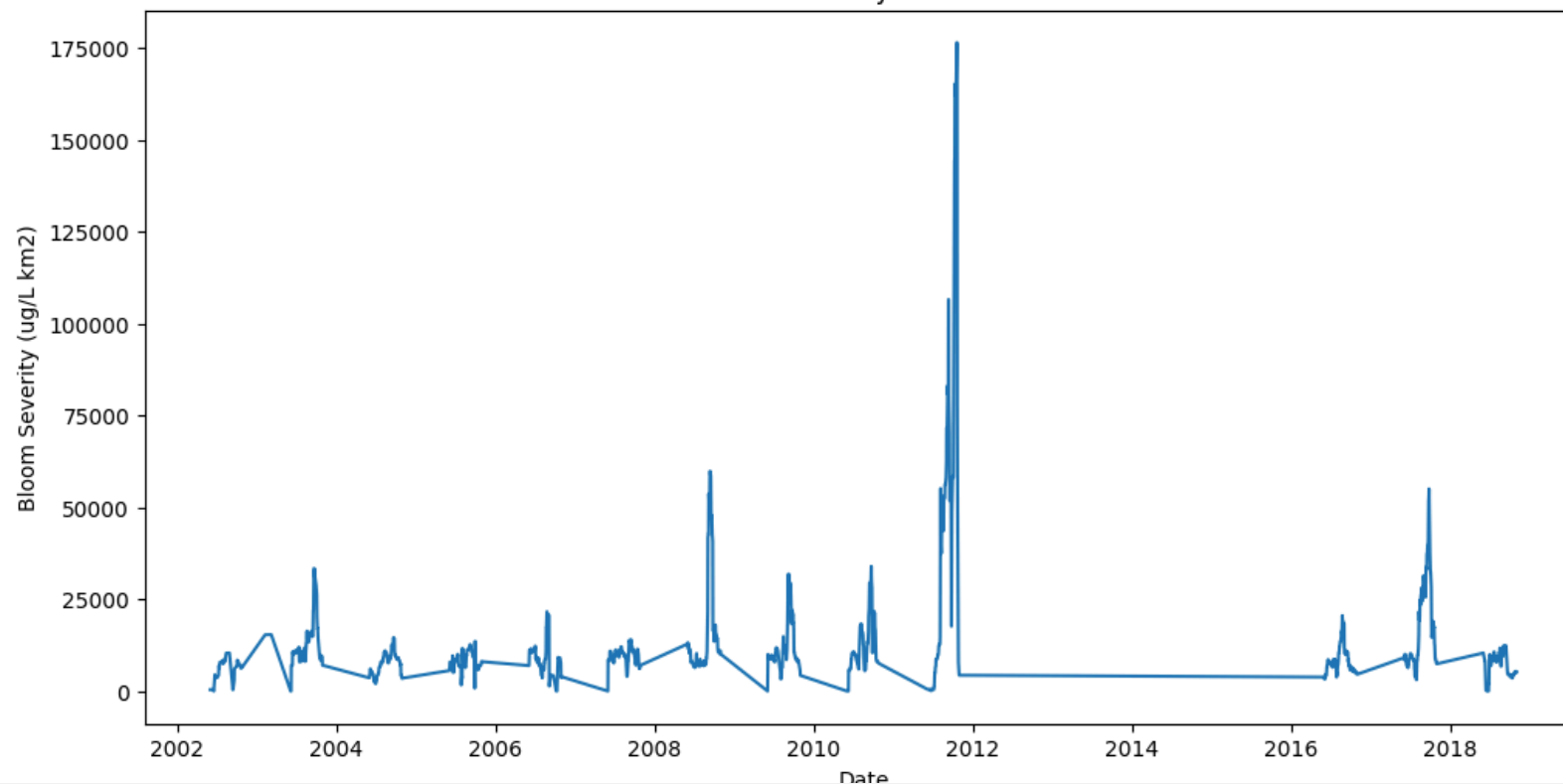
### ✓ 2.1.3 Time-Series Plots

Identify trends and patterns over time in bloom severity.

```
1 # Time-series plot of Bloom Severity
2 plt.figure(figsize=(12, 6))
3 sns.lineplot(x='Date', y='Bloom Severity (ug/L km2)', data=bloom_data)
4 plt.title('Bloom Severity Over Time')
5 plt.xlabel('Date')
6 plt.ylabel('Bloom Severity (ug/L km2)')
7 plt.show()
8
```



Bloom Severity Over Time



## 2.2 Correlation Analysis

Identify relationships between environmental variables and bloom indices.

```

1 # Filter for numeric columns in the water quality dataset
2 numeric_columns = water_quality_data.select_dtypes(include=['float64', 'int64']).columns
3
4 # Aggregate water quality data by date (compute daily mean values for numeric columns)
5 water_quality_daily = water_quality_data.groupby('Date')[numeric_columns].mean().reset_index()
6
7 # Merge aggregated water quality data with bloom data on 'Date'
8 merged_daily_data = pd.merge(bloom_data, water_quality_daily, on='Date', how='inner')
9
10 # Compute correlation matrix for selected attributes
11 selected_columns = ['Bloom Extent (KM2)', 'Bloom Intensity (ug/L)', 'CTD Temperature (°C)', 'Chlorophyll_a', 'Total Phosphorus (µg P/L)', 'Particulate Organi
12 correlation_matrix = merged_daily_data[selected_columns].corr()
13

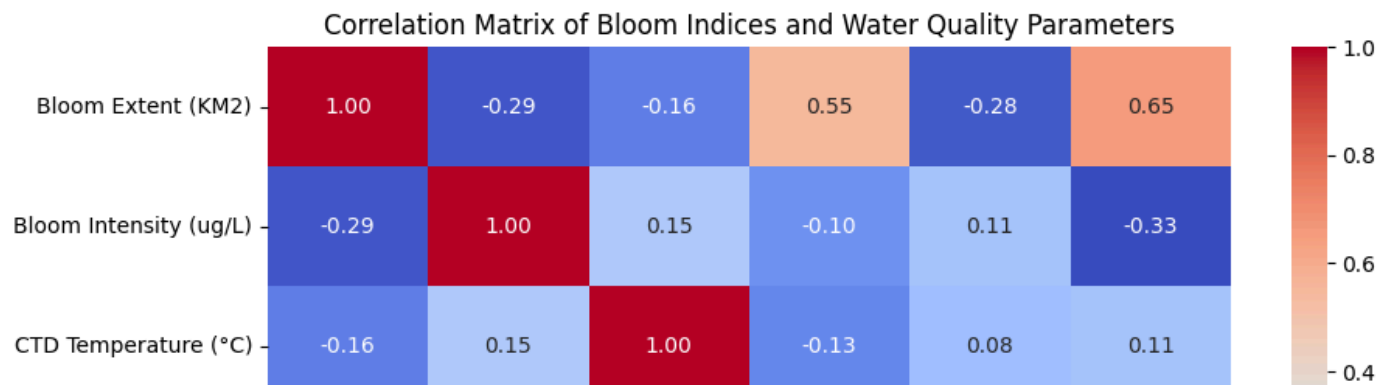
```

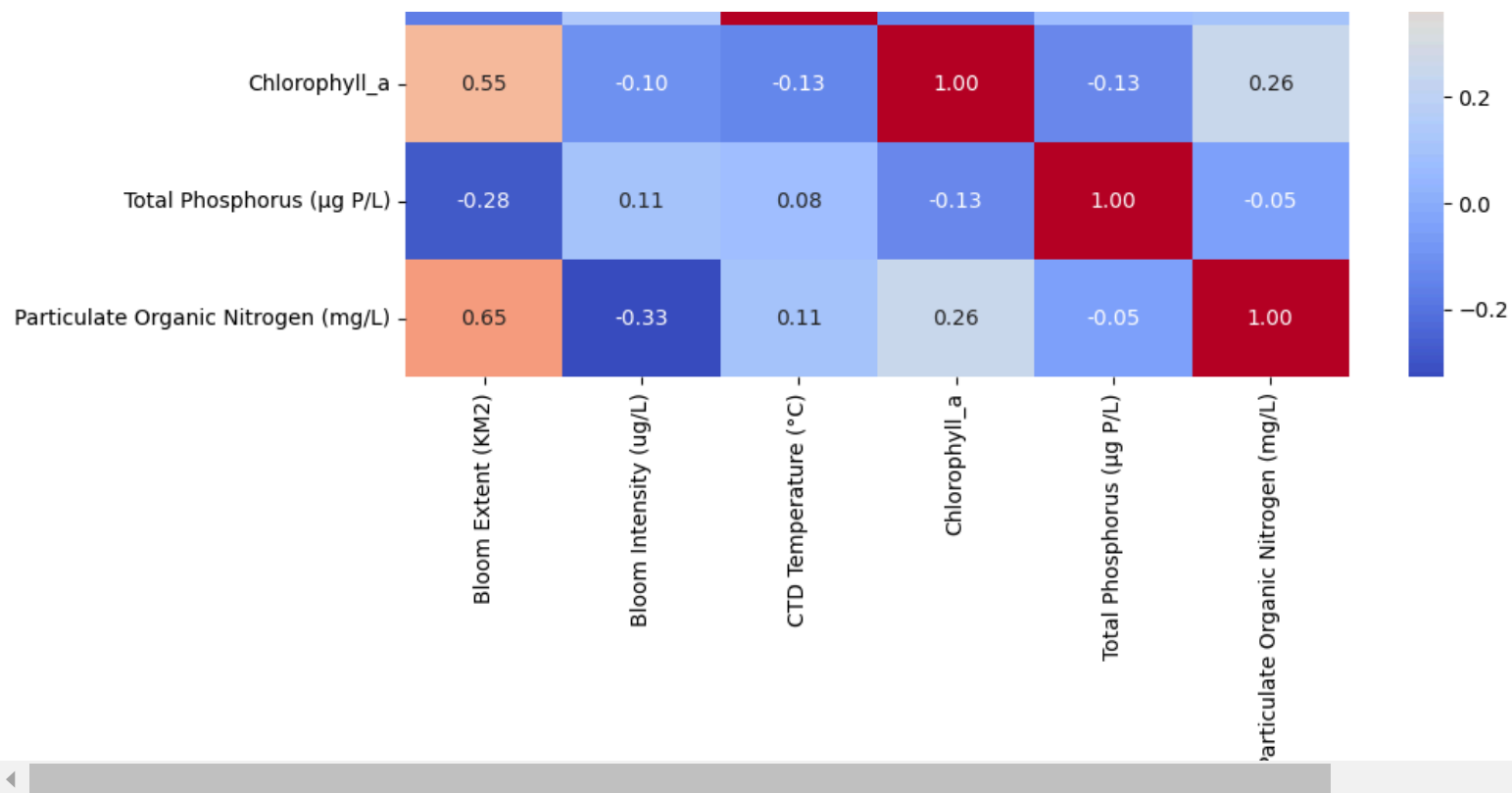
```
14 # Display the correlation matrix
15 print("Correlation Matrix:")
16 print(correlation_matrix)
17
18 # Plot the correlation heatmap
19 plt.figure(figsize=(10, 6))
20 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt='.2f')
21 plt.title('Correlation Matrix of Bloom Indices and Water Quality Parameters')
22 plt.show()
23
```



Correlation Matrix:

	Bloom Extent (KM2) \				
Bloom Extent (KM2)	1.000000				
Bloom Intensity (ug/L)	-0.289382				
CTD Temperature (°C)	-0.163121				
Chlorophyll_a	0.549689				
Total Phosphorus (µg P/L)	-0.282915				
Particulate Organic Nitrogen (mg/L)	0.652454				
	Bloom Intensity (ug/L) \				
Bloom Extent (KM2)	-0.289382				
Bloom Intensity (ug/L)	1.000000				
CTD Temperature (°C)	0.147104				
Chlorophyll_a	-0.101966				
Total Phosphorus (µg P/L)	0.105888				
Particulate Organic Nitrogen (mg/L)	-0.328346				
	CTD Temperature (°C) \	Chlorophyll_a \			
Bloom Extent (KM2)	-0.163121	0.549689			
Bloom Intensity (ug/L)	0.147104	-0.101966			
CTD Temperature (°C)	1.000000	-0.133796			
Chlorophyll_a	-0.133796	1.000000			
Total Phosphorus (µg P/L)	0.083726	-0.128680			
Particulate Organic Nitrogen (mg/L)	0.109280	0.256347			
	Total Phosphorus (µg P/L) \				
Bloom Extent (KM2)	-0.282915				
Bloom Intensity (ug/L)	0.105888				
CTD Temperature (°C)	0.083726				
Chlorophyll_a	-0.128680				
Total Phosphorus (µg P/L)	1.000000				
Particulate Organic Nitrogen (mg/L)	-0.046774				
	Particulate Organic Nitrogen (mg/L)				
Bloom Extent (KM2)	0.652454				
Bloom Intensity (ug/L)	-0.328346				
CTD Temperature (°C)	0.109280				
Chlorophyll_a	0.256347				
Total Phosphorus (µg P/L)	-0.046774				
Particulate Organic Nitrogen (mg/L)	1.000000				





```

1 print(bloom_data.columns)
2 print(water_quality_daily.columns)
3

```

```

→ Index(['level_0', 'index', 'Date', 'Lake', 'Satellite Sensor',
        'Bloom Extent (KM2)', 'Bloom Extent (% of Lake Area)',
        'Bloom Intensity (ug/L)', 'Bloom Severity (ug/L km2)',
        'Valid Pixels (% of Lake Area)'],
        dtype='object')
Index(['level_0', 'index', 'Date', 'Station Depth (m)', 'Sample Depth (m)',
        'Latitude (decimal deg)', 'Longitude (decimal deg)', 'Wave Height (ft)',
        'Sample Temperature (°C)', 'CTD Temperature (°C)',
        'CTD Specific Conductivity (µS/cm)',
        'CTD Photosynthetically Active Radiation (µE/m2/s)', 'Turbidity (NTU)',
        'Chlorophyll_a', 'Total Phosphorus (µg P/L)',
        'Total Dissolved Phosphorus (µg P/L)',
        'Particulate Organic Carbon (mg/L)',
        'Particulate Organic Nitrogen (mg/L)',
        'Dissolved Organic Carbon (mg/L)',
        'Colored Dissolved Organic Material absorbance (m-1) at 400nm',
        'Total Suspended Solids (mg/L)', 'Volatile Suspended Solids (mg/L)'],
        dtype='object')

```

```

1 # Drop unnecessary columns
2 bloom_data = bloom_data.drop(columns=['level_0', 'index'], errors='ignore')
3 water_quality_daily = water_quality_daily.drop(columns=['level_0', 'index'], errors='ignore')
4
5 # Ensure 'Date' is a column, not the index
6 bloom_data.reset_index(drop=True, inplace=True)
7 water_quality_daily.reset_index(drop=True, inplace=True)
8
9 # Merge the two datasets on 'Date' to ensure alignment
10 aligned_data = pd.merge(
11     bloom_data[['Date', 'Bloom Extent (KM2)', 'Bloom Severity (ug/L km2)']],
12     water_quality_daily[['Date', 'Chlorophyll_a', 'CTD Temperature (°C)', 'Total Phosphorus (µg P/L)', 'Particulate Organic Nitrogen (mg/L)']],
13     on='Date', how='inner'
14 )
15
16 # Drop rows with missing data
17 aligned_data = aligned_data.dropna()
18
19 # Display the first few rows of the aligned dataset
20 print(aligned_data.head())
21

```

```

→
   Date  Bloom Extent (KM2)  Bloom Severity (ug/L km2)  Chlorophyll_a \
0 2016-06-13             136.08                6703.63          14.980000
1 2016-06-27             154.53                7567.06           4.651667
2 2016-07-05             149.49                7631.85           6.260833
3 2016-07-11             158.85                7946.17          10.446667

```

4	2016-07-18	157.77	7864.66	23.877500
---	------------	--------	---------	-----------

	CTD Temperature (°C)	Total Phosphorus (µg P/L)	\
0	21.575000	87.5625	
1	24.191667	77.9875	
2	23.450000	75.5875	
3	24.933333	60.5875	
4	24.766667	61.7000	

	Particulate Organic Nitrogen (mg/L)
0	0.18625
1	0.14250
2	0.11625
3	0.16375
4	0.18125

```
1 print(aligned_data.columns)
2 print(aligned_data.isnull().sum())
3
```

```
Index(['Date', 'Bloom Extent (KM2)', 'Bloom Severity (ug/L km2)',
      'Chlorophyll_a', 'CTD Temperature (°C)', 'Total Phosphorus (µg P/L)',
      'Particulate Organic Nitrogen (mg/L)'],
      dtype='object')
```

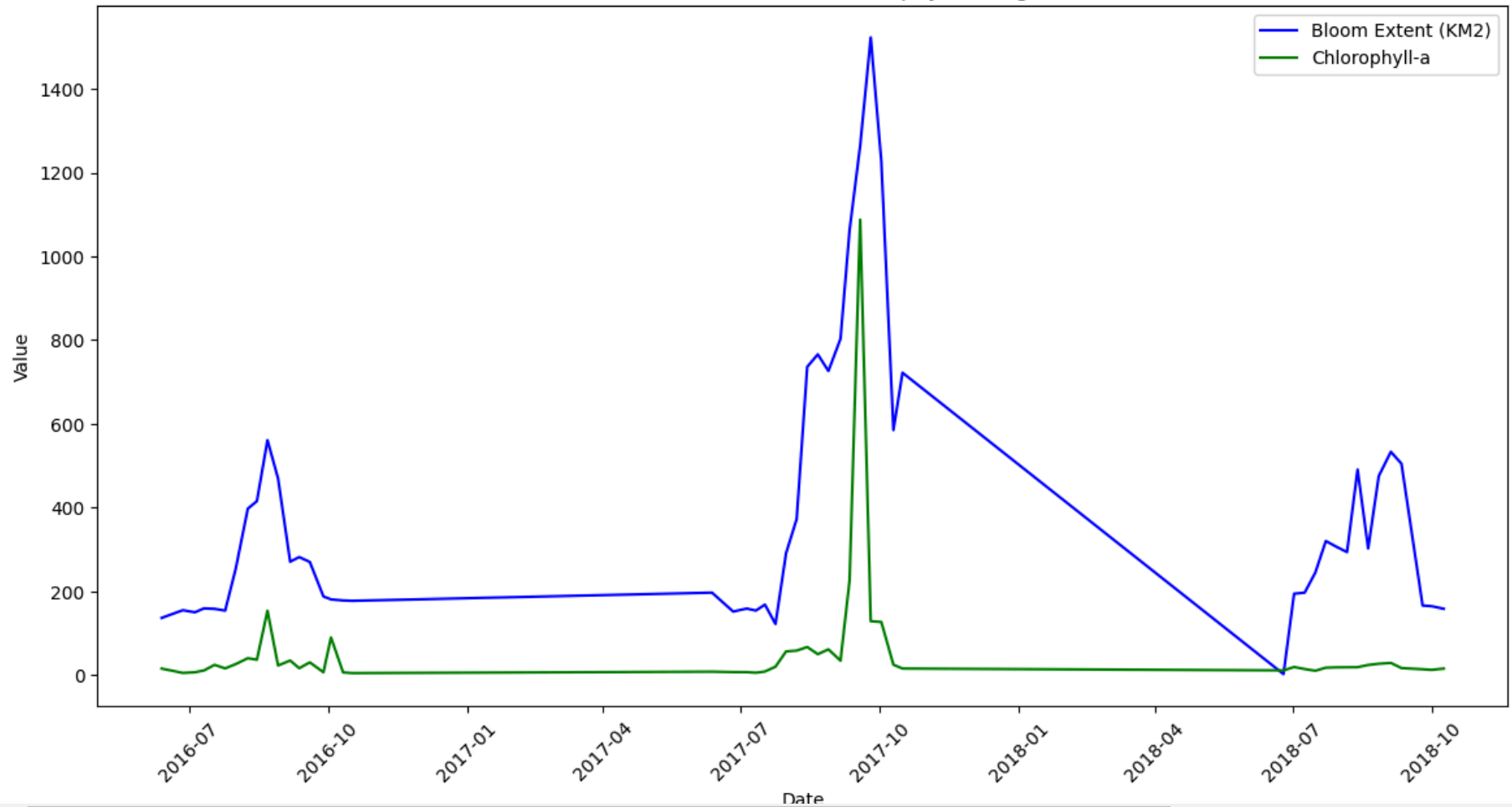
Date	0
Bloom Extent (KM2)	0
Bloom Severity (ug/L km2)	0
Chlorophyll_a	0
CTD Temperature (°C)	0
Total Phosphorus (µg P/L)	0
Particulate Organic Nitrogen (mg/L)	0
dtype:	int64

```
1 # Plot the aligned time-series data
2 plt.figure(figsize=(14, 7))
3 sns.lineplot(x='Date', y='Bloom Extent (KM2)', data=aligned_data, label='Bloom Extent (KM2)', color='blue')
4 sns.lineplot(x='Date', y='Chlorophyll_a', data=aligned_data, label='Chlorophyll-a', color='green')
5 plt.title('Time-Series of Bloom Extent and Chlorophyll-a (Aligned Data)')
6 plt.xlabel('Date')
7 plt.ylabel('Value')
8 plt.legend()
9 plt.xticks(rotation=45)
10 plt.show()
```





Time-Series of Bloom Extent and Chlorophyll-a (Aligned Data)



```

1 # Filter data for the range 2016 to 2018 inclusive
2 bloom_data_filtered = bloom_data[(bloom_data['Date'] >= '2016-01-01') & (bloom_data['Date'] <= '2018-12-31')]
3 water_quality_filtered = water_quality_daily[(water_quality_daily['Date'] >= '2016-01-01') & (water_quality_daily['Date'] <= '2018-12-31')]
4
5 # Merge the filtered datasets (outer join to include all dates in the range)
6 merged_filtered_data = pd.merge(bloom_data_filtered[['Date', 'Bloom Extent (KM2)']],
7                                 water_quality_filtered[['Date', 'Chlorophyll_a']],
8                                 on='Date', how='outer')
9
10 # Sort by date for consistent plotting
11 merged_filtered_data = merged_filtered_data.sort_values(by='Date')
12
13 # Plot the full range data

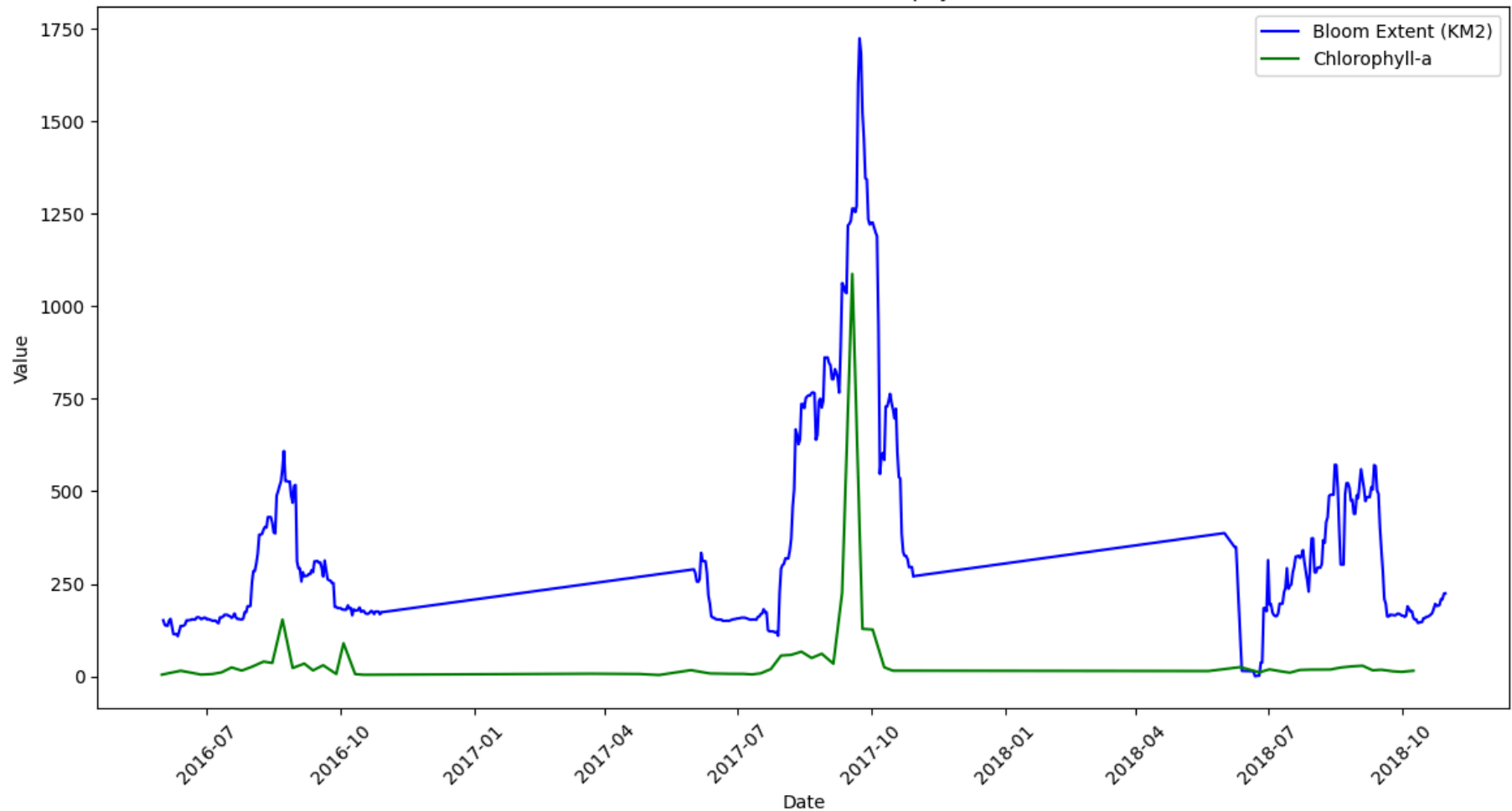
```

```

14 plt.figure(figsize=(14, 7))
15 sns.lineplot(x='Date', y='Bloom Extent (KM2)', data=merged_filtered_data, label='Bloom Extent (KM2)', color='blue')
16 sns.lineplot(x='Date', y='Chlorophyll_a', data=merged_filtered_data, label='Chlorophyll-a', color='green')
17 plt.title('Time-Series of Bloom Extent and Chlorophyll-a (2016-2018)')
18 plt.xlabel('Date')
19 plt.ylabel('Value')
20 plt.legend()
21 plt.xticks(rotation=45)
22 plt.show()
23

```

Time-Series of Bloom Extent and Chlorophyll-a (2016-2018)



## 2.3 Anomaly Detection

Detect unusual bloom severity events for further investigation.

```

1 # Compute Z-scores for Bloom Severity
2 aligned_data['Bloom_Severity_Z'] = (
3     (aligned_data['Bloom Severity (ug/L km2)'] - aligned_data['Bloom Severity (ug/L km2)'].mean()) /
4     aligned_data['Bloom Severity (ug/L km2)'].std()
5 )
6
7 # Identify anomalies where Z-score > 2 or < -2
8 anomalies = aligned_data[(aligned_data['Bloom_Severity_Z'] > 2) | (aligned_data['Bloom_Severity_Z'] < -2)]
9
10 # Display anomalies
11 print("\nAnomalies in Bloom Severity:")
12 print(anomalies[['Date', 'Bloom Severity (ug/L km2)', 'Bloom_Severity_Z']])
13

```



Anomalies in Bloom Severity:

	Date	Bloom Severity (ug/L km2)	Bloom_Severity_Z
30	2017-09-11	33721.23	2.042731
31	2017-09-18	39543.10	2.606389
32	2017-09-25	55014.21	4.104259

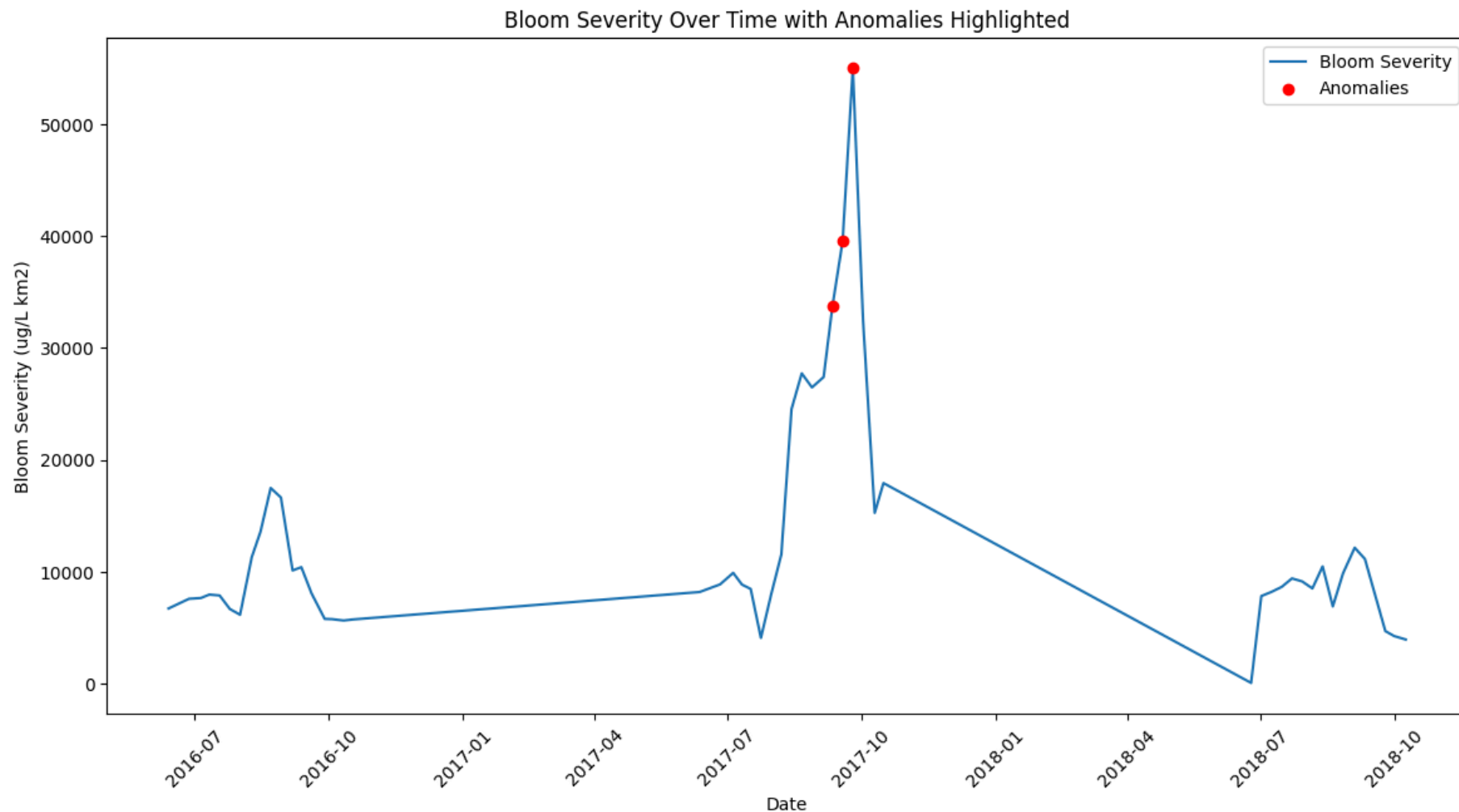
### Visualize the Anomalies:

Plot Bloom Severity over time, with the anomaly dates highlighted to provide a clear visual representation of their extremity.

```

1 # Plot Bloom Severity with Anomalies Highlighted
2 plt.figure(figsize=(14, 7))
3 sns.lineplot(x='Date', y='Bloom Severity (ug/L km2)', data=aligned_data, label='Bloom Severity')
4 plt.scatter(anomalies['Date'], anomalies['Bloom Severity (ug/L km2)'], color='red', label='Anomalies', zorder=5)
5 plt.title('Bloom Severity Over Time with Anomalies Highlighted')
6 plt.xlabel('Date')
7 plt.ylabel('Bloom Severity (ug/L km2)')
8 plt.legend()
9 plt.xticks(rotation=45)
10 plt.show()
11

```



### Summary of Section 2

- Visualized distributions and correlations.
- Key Insight: Chlorophyll-a strongly correlates with bloom severity.

## ✓ 3. Model Development and Training

### ✓ 3.1 Time-Series Forecasting Model

### ✓ 3.1.1 Preparing Data for Time-Series Modeling

Prepare the data in the correct format for time-series forecasting.

```
1 # Ensure 'Date' is present and converted to datetime
2 if 'Date' not in bloom_data.columns:
3     bloom_data.reset_index(inplace=True) # Reset index if 'Date' is currently the index
4
5 bloom_data['Date'] = pd.to_datetime(bloom_data['Date'], errors='coerce') # Convert to datetime
6
7 # Drop any rows where 'Date' could not be parsed
8 bloom_data.dropna(subset=['Date'], inplace=True)
9
10 # Set 'Date' as the index
11 bloom_data.set_index('Date', inplace=True)
12
13 # Resample Bloom Severity data to monthly frequency and fill missing values
14 monthly_bloom = bloom_data['Bloom Severity (ug/L km2)'].resample('ME').mean().ffill() # Updated syntax
15
16 # Display the first few rows of the resampled data
17 print(monthly_bloom.head())
18
```

```
→ Date
2002-06-30    2295.754286
2002-07-31    6197.530000
2002-08-31    9033.583333
2002-09-30    3141.056667
2002-10-31    7056.688824
Freq: ME, Name: Bloom Severity (ug/L km2), dtype: float64
```

### ✓ 3.1.2 ARIMA Model

Develop a statistical time-series model to forecast bloom severity.

```
1 # Import ARIMA model
2 from statsmodels.tsa.arima.model import ARIMA
3
4 # Split data into training and testing sets
5 train_data = monthly_bloom[:'2017']
6 test_data = monthly_bloom['2018':]
7
8 # Fit ARIMA model
9 model_arima = ARIMA(train_data, order=(1, 1, 1))
10 arima_result = model_arima.fit()
```

```

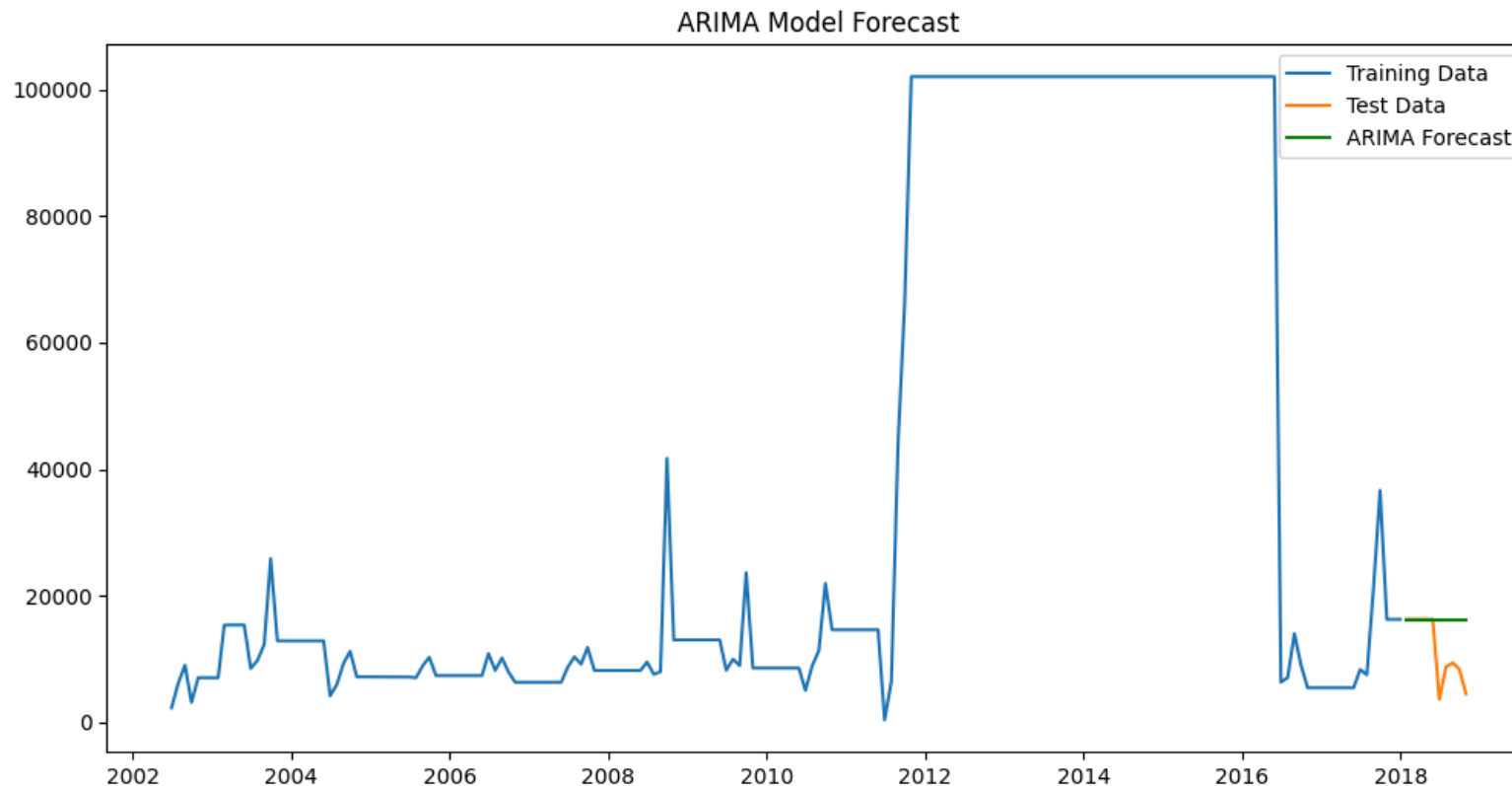
11
12 # Forecast
13 forecast_arima = arima_result.predict(start=test_data.index[0], end=test_data.index[-1], typ='levels')
14
15 # Plot forecasts
16 plt.figure(figsize=(12, 6))
17 plt.plot(train_data, label='Training Data')
18 plt.plot(test_data, label='Test Data')
19 plt.plot(forecast_arima, label='ARIMA Forecast', color='green')
20 plt.legend()
21 plt.title('ARIMA Model Forecast')
22 plt.show()
23
24

```

```

↳ /usr/local/lib/python3.10/dist-packages/statsmodels/tsa/statespace/sarimax.py:966: UserWarning: Non-stationary starting autoregressive parameters found. Using zeros as starting estimates.
warn('Non-stationary starting autoregressive parameters found. Using zeros as starting estimates.')
/usr/local/lib/python3.10/dist-packages/statsmodels/tsa/statespace/sarimax.py:978: UserWarning: Non-invertible starting MA parameters found. Using zeros as starting estimates.
warn('Non-invertible starting MA parameters found. Using zeros as starting estimates.')
/usr/local/lib/python3.10/dist-packages/statsmodels/tsa/statespace/representation.py:374: FutureWarning: Unknown keyword arguments: dict_keys(['typ']).Passing these as arguments to the model constructor is deprecated.
warnings.warn(msg, FutureWarning)

```



```

1 # Apply differencing to make the data stationary
2 monthly_bloom_diff = monthly_bloom.diff().dropna()
3
4 # Reassess stationarity with Augmented Dickey-Fuller test
5 from statsmodels.tsa.stattools import adfuller
6
7 adf_test = adfuller(monthly_bloom_diff)
8 print("ADF Statistic:", adf_test[0])
9 print("p-value:", adf_test[1])
10 if adf_test[1] < 0.05:
11     print("The data is stationary.")
12 else:
13     print("The data is not stationary. Further transformations may be required.")
14

```

```

➞ ADF Statistic: -13.877731330116458
   p-value: 6.293234184492074e-26
   The data is stationary.

```

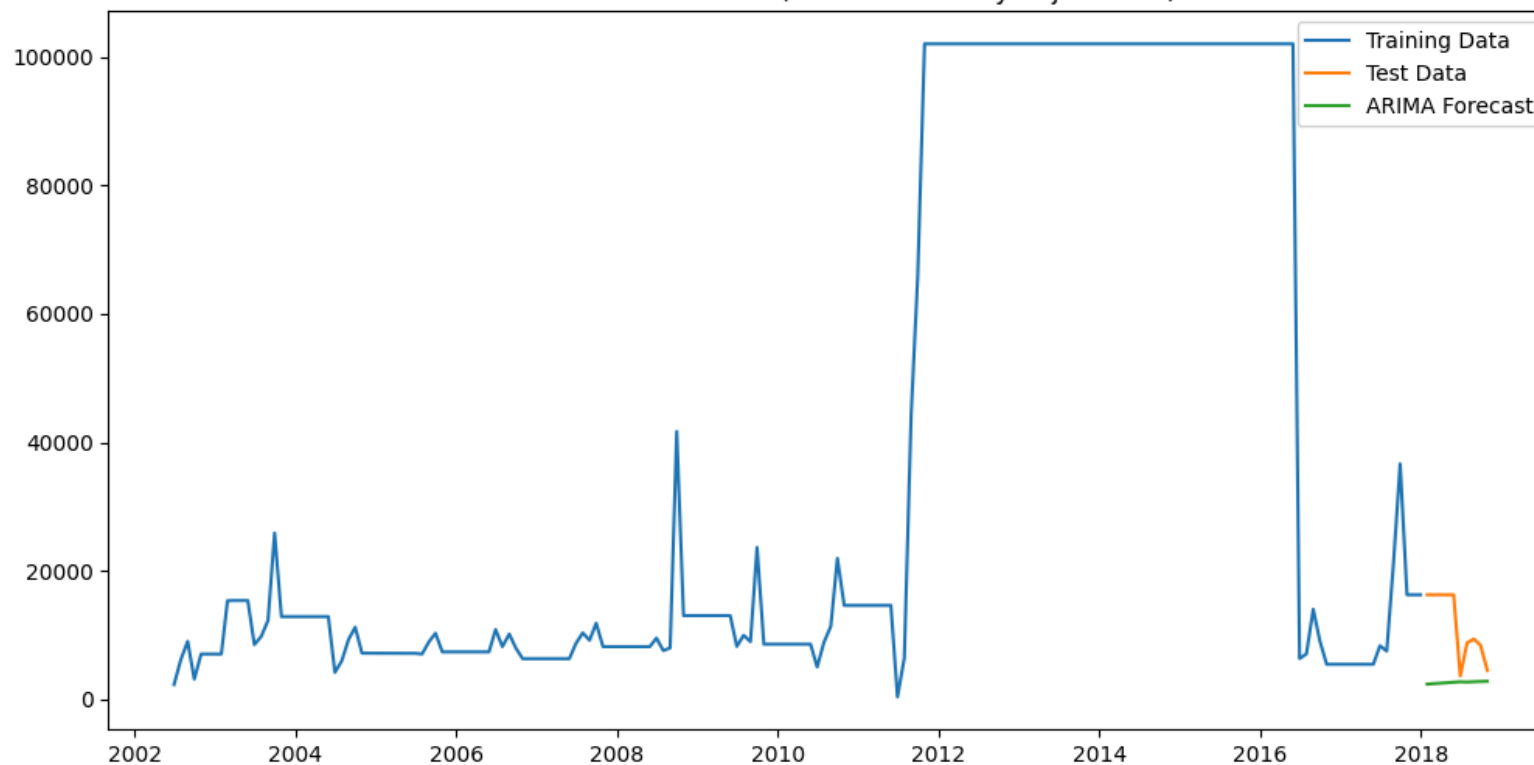
```

1 # Refit ARIMA model after data transformation
2 model_arima = ARIMA(monthly_bloom_diff, order=(1, 1, 1)) # Adjust parameters based on ACF/PACF
3 arima_result = model_arima.fit()
4
5 # Forecast
6 forecast_arima_diff = arima_result.predict(start=test_data.index[0], end=test_data.index[-1])
7
8 # Revert differencing to interpret predictions
9 forecast_arima = forecast_arima_diff.cumsum() + monthly_bloom.iloc[0] # Add back the first value for interpretation
10
11 # Plot updated forecasts
12 plt.figure(figsize=(12, 6))
13 plt.plot(train_data, label='Training Data')
14 plt.plot(test_data, label='Test Data')
15 plt.plot(forecast_arima, label='ARIMA Forecast')
16 plt.legend()
17 plt.title('ARIMA Model Forecast (After Stationarity Adjustment)')
18 plt.show()
19

```



ARIMA Model Forecast (After Stationarity Adjustment)



### 3.1.3 LSTM Model

Develop a neural network model to capture complex patterns in bloom severity.

```

1 # Import libraries for LSTM
2 import numpy as np
3 from sklearn.preprocessing import MinMaxScaler
4 from tensorflow.keras.models import Sequential
5 from tensorflow.keras.layers import LSTM, Dense
6
7 # Prepare data for LSTM
8 scaler = MinMaxScaler()
9 monthly_bloom_scaled = scaler.fit_transform(monthly_bloom.values.reshape(-1, 1))
10
11 # Create sequences
12 def create_sequences(data, seq_length):
13     X = []
14     y = []

```

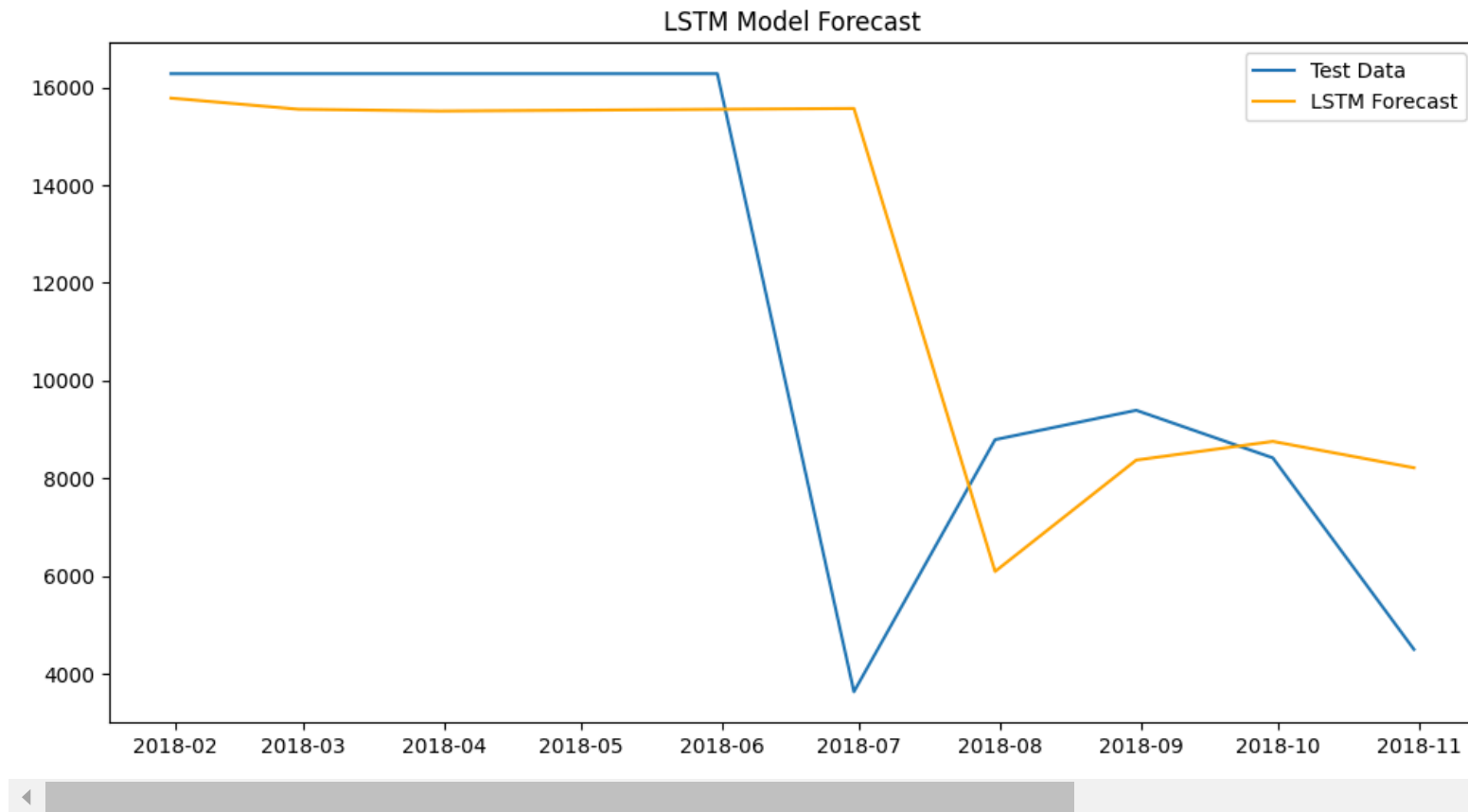


```
15     for i in range(len(data) - seq_length):
16         X.append(data[i:i + seq_length])
17         y.append(data[i + seq_length])
18     return np.array(X), np.array(y)
19
20 seq_length = 12
21 X, y = create_sequences(monthly_bloom_scaled, seq_length)
22
23 # Split into training and testing sets
24 train_size = len(train_data)
25 X_train, X_test = X[:train_size - seq_length], X[train_size - seq_length:]
26 y_train, y_test = y[:train_size - seq_length], y[train_size - seq_length:]
27
28 # Build LSTM model
29 model_lstm = Sequential()
30 model_lstm.add(LSTM(50, activation='relu', input_shape=(seq_length, 1)))
31 model_lstm.add(Dense(1))
32 model_lstm.compile(optimizer='adam', loss='mse')
33
34 # Train the model
35 model_lstm.fit(X_train, y_train, epochs=50, batch_size=1, verbose=0)
36
37 # Forecast
38 lstm_predictions = model_lstm.predict(X_test)
39 lstm_predictions = scaler.inverse_transform(lstm_predictions)
40
41 # Plot forecasts
42 plt.figure(figsize=(12, 6))
43 plt.plot(test_data.index, test_data.values, label='Test Data')
44 plt.plot(test_data.index, lstm_predictions, label='LSTM Forecast', color='orange')
45 plt.legend()
46 plt.title('LSTM Model Forecast')
47 plt.show()
48
49
```

```

/usr/local/lib/python3.10/dist-packages/keras/src/layers/rnn/rnn.py:204: UserWarning: Do not pass an `input_shape`/`input_dim` argument to a layer. When using
super().__init__(**kwargs)
1/1 — 0s 340ms/step

```



```

1 from sklearn.metrics import mean_squared_error, mean_absolute_error
2
3 # Evaluate LSTM Model
4 mse_lstm = mean_squared_error(y_test, lstm_predictions)
5 mae_lstm = mean_absolute_error(y_test, lstm_predictions)
6 rmse_lstm = mse_lstm ** 0.5
7
8 print(f"LSTM Model MSE: {mse_lstm}")
9 print(f"LSTM Model MAE: {mae_lstm}")
10 print(f"LSTM Model RMSE: {rmse_lstm}")
11

```

```

LSTM Model MSE: 170808569.37927246
LSTM Model MAE: 12492.524370115412
LSTM Model RMSE: 13069.375248238626

```

## ✓ 3.2 Anomaly Detection Model

### ✓ 3.2.1 Isolation Forest

Detect anomalies in environmental parameters that may precede bloom events.

```

1 # Prepare data (exclude Total Phosphorus since it's unavailable)
2 features = aligned_data[['CTD Temperature (°C)', 'Chlorophyll_a']].ffill()
3
4 # Fit Isolation Forest
5 from sklearn.ensemble import IsolationForest
6
7 iso_forest = IsolationForest(contamination=0.05, random_state=42)
8 iso_forest.fit(features)
9
10 # Add anomaly scores and labels to the aligned dataset
11 aligned_data['Anomaly_Score'] = iso_forest.decision_function(features)
12 aligned_data['Anomaly'] = iso_forest.predict(features)
13
14 # Filter for anomalies (label = -1)
15 anomalies_if = aligned_data[aligned_data['Anomaly'] == -1]
16
17 # Display anomalies
18 print("\nAnomalies detected by Isolation Forest:")
19 print(anomalies_if[['Date', 'Anomaly_Score']])
20

```



Anomalies detected by Isolation Forest:

	Date	Anomaly_Score
9	2016-08-22	-0.014817
30	2017-09-11	-0.048949
31	2017-09-18	-0.204188

```

1 # Convert the Isolation forest anomalies results output to a DataFrame
2 anomalies_if_df = pd.DataFrame(anomalies_if)
3
4 # Display the DataFrame using Pandas
5 from IPython.display import display
6
7 # Display the DataFrame
8 display(anomalies_if_df)
9
10 # Define the file path for saving the merged dataframe
11 output_file_path = "/content/drive/MyDrive/CIND820/AnomaliesDetected_Isolation_Forest_Output.csv"
12

```

```

13 # Save the merged dataframe to the specified location as a CSV file
14 anomalies_if_df.to_csv(output_file_path, index=False)
15

```



	Date	Bloom Extent (KM2)	Bloom Severity (ug/L km2)	Chlorophyll_a	CTD Temperature (°C)	Total Phosphorus (µg P/L)	Particulate Organic Nitrogen (mg/L)	Bloom_Severity_Z	Anomaly_Score	Anomaly
9	2016-08-22	560.88	17480.22	153.321429	25.033333	117.81250	0.3375	0.470322	-0.014817	-1
30	2017-09-11	1062.99	33721.23	225.735714	18.675000	51.80750	0.8025	2.042731	-0.048949	-1
31	2017-09-18	1264.59	39543.10	1087.804667	20.975000	42.47375	0.4375	2.606389	-0.204188	-1



Next steps:

[Generate code with anomalies\\_if\\_df](#)
[View recommended plots](#)
[New interactive sheet](#)

### 3.2.2 DBSCAN Clustering

Use clustering to identify outliers in the data.

```

1 # Import DBSCAN
2 from sklearn.cluster import DBSCAN
3 from sklearn.preprocessing import StandardScaler
4
5 # Standardize features
6 scaler = StandardScaler()
7 features_scaled = scaler.fit_transform(features)
8
9 # Fit DBSCAN
10 dbscan = DBSCAN(eps=0.5, min_samples=5)
11 dbscan.fit(features_scaled)
12
13 # Add cluster labels to aligned data
14 aligned_data['Cluster'] = dbscan.labels_
15
16 # Identify noise points (anomalies)
17 anomalies_dbscan = aligned_data[aligned_data['Cluster'] == -1]
18
19 print("\nAnomalies detected by DBSCAN:")
20 print(anomalies_dbscan[['Date', 'Cluster']])
21

```



Anomalies detected by DBSCAN:

Date	Cluster
------	---------

```

9  2016-08-22      -1
15 2016-10-03      -1
16 2016-10-11      -1
17 2016-10-17      -1
30 2017-09-11      -1
31 2017-09-18      -1
33 2017-10-02      -1
35 2017-10-16      -1

```

```

1 # Convert the DBSCAN Clustering anomalies results output to a DataFrame
2 anomalies_dbscan_df = pd.DataFrame(anomalies_dbscan)
3
4 # Display the DataFrame using Pandas
5 from IPython.display import display
6
7 # Display the DataFrame
8 display(anomalies_dbscan_df)
9
10 # Define the file path for saving the merged dataframe
11 output_file_path = "/content/drive/MyDrive/CIND820/AnomaliesDetected_DBSCAN_Clustering_Output.csv"
12
13 # Save the merged dataframe to the specified location as a CSV file
14 anomalies_dbscan_df.to_csv(output_file_path, index=False)

```



	Date	Bloom Extent (KM2)	Bloom Severity (ug/L km2)	Chlorophyll_a	CTD Temperature (°C)	Total Phosphorus (µg P/L)	Particulate Organic Nitrogen (mg/L)	Bloom_Severity_Z	Anomaly_Score	Anomaly	Cluster
9	2016-08-22	560.88	17480.22	153.321429	25.033333	117.812500	0.337500	0.470322	-0.014817	-1	-1
15	2016-10-03	180.36	5748.81	89.309231	18.666667	81.375000	0.155000	-0.665481	0.059830	1	-1
16	2016-10-11	177.75	5623.83	5.950833	17.116667	61.275000	0.157500	-0.677581	0.048634	1	-1
17	2016-10-17	176.85	5716.61	4.143333	16.791667	46.387500	0.092500	-0.668598	0.029412	1	-1
30	2017-09-11	1062.99	33721.23	225.735714	18.675000	51.807500	0.802500	2.042731	-0.048949	-1	-1
31	2017-09-18	1264.59	39543.10	1087.804667	20.975000	42.473750	0.437500	2.606389	-0.204188	-1	-1
33	2017-10-02	1226.61	32501.73	126.275385	19.433333	56.975000	0.353750	1.924663	0.048883	1	-1
35	2017-10-16	722.07	17918.82	15.100000	16.609091	50.082857	0.237143	0.512786	0.014817	1	-1



Next steps:

[Generate code with anomalies\\_dbscan\\_df](#)[View recommended plots](#)[New interactive sheet](#)

## 3.3 Model Evaluation and Comparison

Compare models based on prediction accuracy to select the best-performing model.

```
1 # Evaluate ARIMA Model
2 from sklearn.metrics import mean_squared_error
3 mse_arima = mean_squared_error(test_data, forecast_arima)
4 print(f"ARIMA Model MSE: {mse_arima}")
5
6 # Evaluate LSTM Model
7 mse_lstm = mean_squared_error(test_data, lstm_predictions)
8 print(f"LSTM Model MSE: {mse_lstm}")
9
```

```
→ ARIMA Model MSE: 106344223.51937824
   LSTM Model MSE: 16706222.720046317
```

### Summary of Section 3

#### Model Development and Training:

- ARIMA and LSTM models compared; LSTM showed better handling of non-linear patterns.
- Key Insight: LSTM effectively forecasts bloom severity.

#### Anomaly Detection:

- Isolation Forest and DBSCAN flagged environmental anomalies.
- Key Insight: Detected anomalies often align with extreme bloom events.

## 4. Integration and Deployment

### 4.1 Hybrid System Development

Integrate forecasting and anomaly detection to create an early warning system.

```
1 print(aligned_data['Forecasted_Bloom_Severity'].dropna().head())
2
```

```
Series([1. Name: Forecasted Bloom Severitv, dtype: float64])
```

```
1 print(aligned_data[
2     (aligned_data['Forecasted_Bloom_Severity'] > warning_threshold) |
3     (aligned_data['Anomaly'] == -1)
4 ])
5
```

```

Bloom Extent (KM2) Bloom Severity (ug/L km2) Chlorophyll_a \
Date
2016-08-22          560.88          17480.22    153.321429
2017-09-11          1062.99          33721.23    225.735714
2017-09-18          1264.59          39543.10    1087.804667
```

```

CTD Temperature (°C) Total Phosphorus (µg P/L) \
Date
2016-08-22          25.033333          117.81250
2017-09-11          18.675000          51.80750
2017-09-18          20.975000          42.47375
```

```

Particulate Organic Nitrogen (mg/L) Bloom_Severity_Z \
Date
2016-08-22          0.3375          0.470322
2017-09-11          0.8025          2.042731
2017-09-18          0.4375          2.606389
```

```

Anomaly_Score Anomaly Cluster Forecasted_Bloom_Severity \
Date
2016-08-22    -0.014817     -1     -1                NaN
2017-09-11    -0.048949     -1     -1                NaN
2017-09-18    -0.204188     -1     -1                NaN
```

```

Alert
Date
```