

✓ OmniXAI in a ML workflow Using Shark Research Data (August 2017-2022)

This tutorial shows how to apply XAI in different stages in a standard ML workflow. The OmniXai library is used in the example. The goal is to use prepared data from transmitters and receivers off of the Coast of Cape Cod to predict shark presence (0 or 1) in the month of August. The data is limited and prepared for use during WiDS Charlotte 2025, and can only be used for educational, not research purposes.

```
1 %%capture
2 import warnings
3 warnings.simplefilter("ignore")
4 !pip install omnixai
5 !pip install kaleido
6 # This default renderer is used for sphinx docs only. Please delete this cell in IPyThor
7 import plotly.io as pio
8 pio.renderers.default = "png"
9 import os
10 import numpy as np
11 import pandas as pd
12 from sklearn.ensemble import RandomForestClassifier
```

The partial dataset used in this example is July 10 to August 10, allowing for a more balanced target class "#sharks".

The dataset in this notebook can be read from: [GitHub](https://github.com/DrPamelaThompson/WiDS-Charlotte-2025/blob/main/XAI_shark_presence_prediction_WiDS_2.ipynb#scroll=10)

```
1 # fetch dataset from github
2 # Convert to DataFrame
3 # URL of the raw dataset without 15% of the data (holdout for prediction and validation
4 url = "https://raw.githubusercontent.com/DrPamelaThompson/WiDS-Charlotte-2025/refs/heads
5 # Read the CSV
6 df = pd.read_csv(url)
7 #drop na not necessary
8 df.info()
```



```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 8887 entries, 0 to 8886
```

```
Data columns (total 22 columns):
```

#	Column	Non-Null Count	Dtype
0	#sharks	8887 non-null	int64
1	dewF	8887 non-null	float64
2	tempF_air	8887 non-null	float64
3	datetime	8887 non-null	object
4	tau_gust_mag	8887 non-null	float64
5	taux_gust	8887 non-null	float64

```

6  tempF_water_cat      8887 non-null  int64
7  stability_value      8887 non-null  float64
8  tempF_water          8887 non-null  float64
9  precip               8887 non-null  float64
10 alpha                8887 non-null  float64
11 Tide_water_level_cat_encoded  8887 non-null  int64
12 unique_shark_count_cat_encoded  8887 non-null  int64
13 Tide_ebb_flood_cat_encoded    8887 non-null  int64
14 press_change_cat_encoded      8887 non-null  int64
15 Light                    8887 non-null  object
16 tempF_Diff_encoded          8887 non-null  int64
17 press                     8887 non-null  float64
18 Moon_Phase_cat             8887 non-null  object
19 Year                       8887 non-null  int64
20 Month                     8887 non-null  int64
21 Day                       8887 non-null  int64

```

dtypes: float64(9), int64(10), object(3)

memory usage: 1.5+ MB

```
1 df.value_counts('#sharks')
```



```

count
#sharks
0      4484
1      4403

```

dtype: int64

```
1 df.tail()
```



```

#sharks  dewF  tempF_air  datetime  tau_gust_mag  taux_gust  tempF_water_
8882      0  65.833333  72.233333  2022-08-10  0.016946  0.013035
      21:30:00+00:00
8883      1  65.366667  72.000000  2022-08-10  0.017251  0.011803
      22:00:00+00:00
8884      1  65.800000  71.400000  2022-08-10  0.020704  0.010926
      22:30:00+00:00
8885      0  66.000000  71.066667  2022-08-10  0.013833  0.008830
      23:00:00+00:00
8886      1  66.766667  71.000000  2022-08-10  0.011247  0.006504
      23:30:00+00:00

```

5 rows × 22 columns

```
1 #check balance of second target class
2 df.value_counts('unique_shark_count_cat_encoded')
```



	count
unique_shark_count_cat_encoded	
0	4484
1	4347
2	56

dtype: int64

July 12 and 13, 2022: Severe weather affected the region, with reports of heavy rainfall and strong winds impacting Cape Cod. Record numbers will be extracted for later investigation.

```
1 # Convert datetime column to datetime format (if not already)
2 df["datetime"] = pd.to_datetime(df["datetime"])
3
4 # Find the index (record number) for July 12, 2022
5 record_index = df[df["datetime"] == "2022-07-12"].index
6
7 # Display the record number(s)
8 print("Record number(s) for July 12, 2022:", record_index.tolist())
```



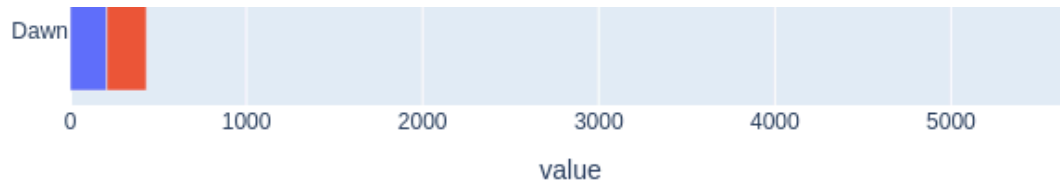
Record number(s) for July 12, 2022: [7529]

The target is #sharks and a second target, not used, is unique_shark_count_cat_encoded

Let's first check if some features are correlated and if there exists data imbalance issues that leads to any potential bias. We can create an `DataAnalyzer` explainer from OmniXai to do this task.

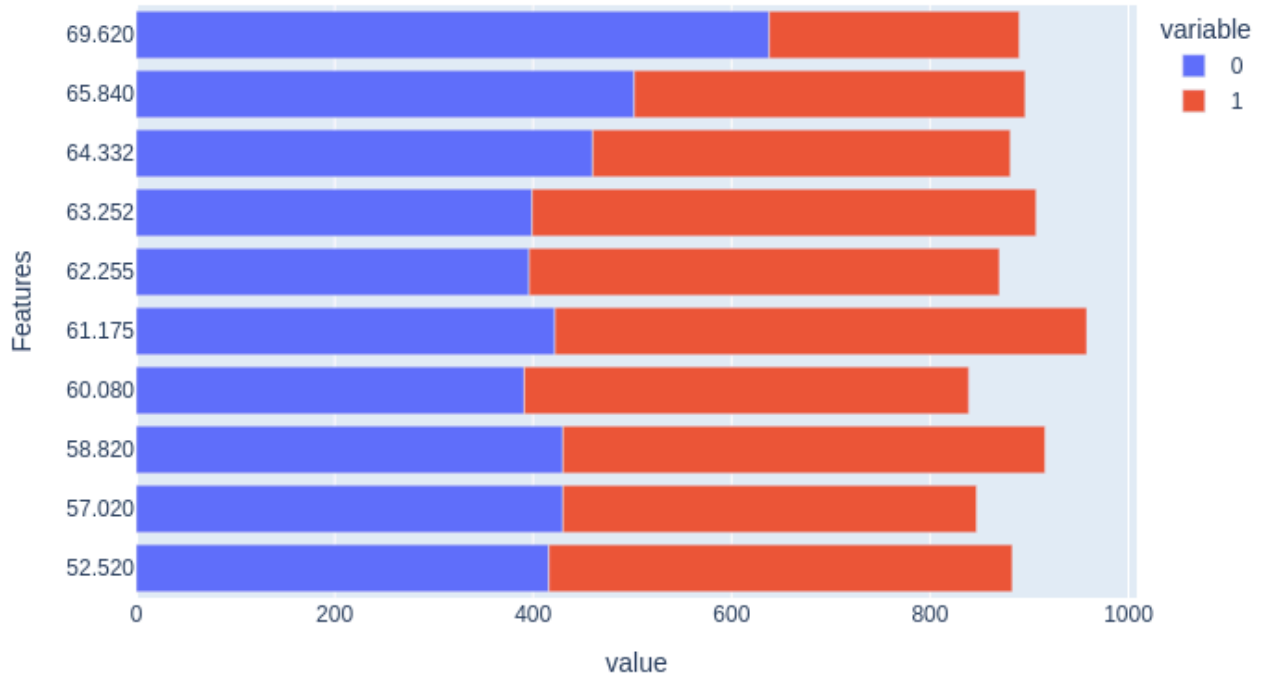
```
1 df1 = df.drop(columns=['datetime', 'unique_shark_count_cat_encoded', 'Year', 'Month', 'Day'])
```

```
1 from omnixai.data.tabular import Tabular
2 from omnixai.explainers.data import DataAnalyzer
3
4 tabular_data = Tabular(
5     df1,
6     categorical_columns=['Light', 'Moon_Phase_cat'],
7     target_column='#sharks'
8 )
9
10 # Initialize a `DataAnalyzer` explainer.
11 # We can choose multiple explainers/analyzers by specifying analyzer names.
12 # In this example, the first explainer is for feature correlation analysis and
13 # the others are for feature imbalance analysis (the same explainer with different param
14 explainer = DataAnalyzer(
15     explainers=["correlation", "imbalance#0", "imbalance#1", "imbalance#2"],
16     mode="classification",
17     data=tabular_data
18 )
19 # Generate explanations by calling `explain_global`.
20 explanations = explainer.explain_global(
21     params={"imbalance#0": {"features": ["Light"]},
22             "imbalance#1": {"features": ["tempF_water"]},
23             "imbalance#2": {"features": ["Moon_Phase_cat"]},
24             "imbalance#3": {"features": ["Light", "tempF_water"]}}
25 )
26
27 print("Correlation:")
28 explanations["correlation"].ipython_plot()
29 print("Imbalance#0: features tempF_water_cat")
30 explanations["imbalance#0"].ipython_plot()
31 print("Imbalance#1: features tempF_water")
32 explanations["imbalance#1"].ipython_plot()
33 print("Imbalance#2: features Light, tempF_water")
34 explanations["imbalance#2"].ipython_plot()
```



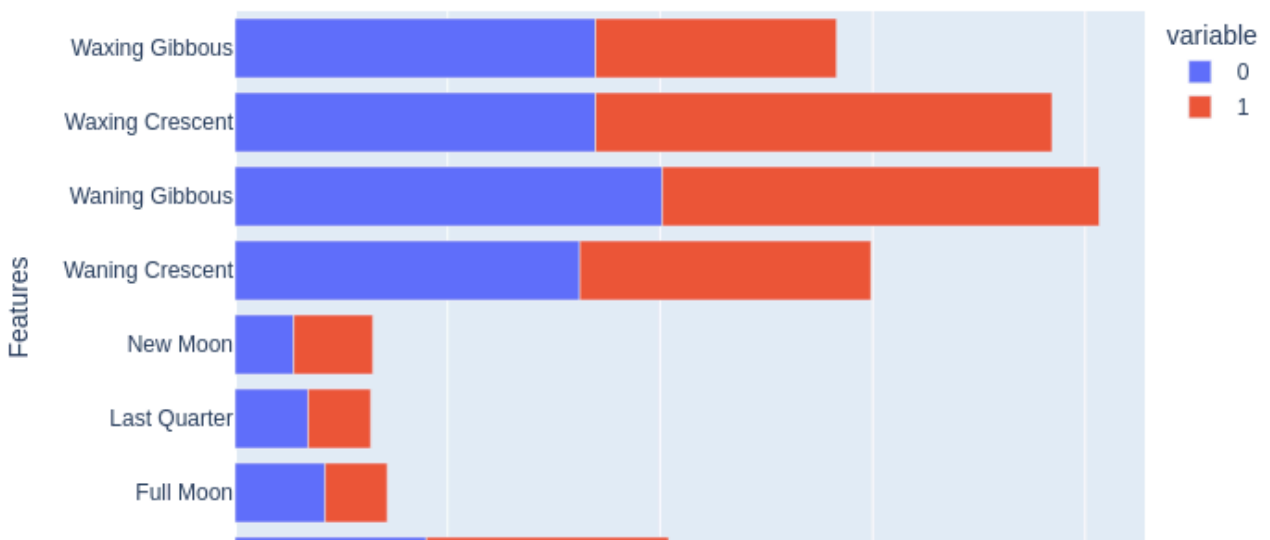
Imbalance#1: features tempF_water

Imbalance Plot



Imbalance#2: features Light, tempF_water

Imbalance Plot





From the correlation plot we can observe that "tempF_water" has strong correlations with "tempF_water_cat", so we may remove one of these features. We will also remove tau_gust based on domain information. From the data imbalance plots we can see that the class labels are relatively balanced in the features.

```
1 df2 = df1.drop(columns=["tempF_water_cat", "tau_gust"])
2 df2.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8887 entries, 0 to 8886
Data columns (total 15 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   #sharks                               8887 non-null   int64
1   dewF                                  8887 non-null   float64
2   tempF_air                             8887 non-null   float64
3   tau_gust_mag                          8887 non-null   float64
4   stability_value                       8887 non-null   float64
5   tempF_water                           8887 non-null   float64
6   precip                                8887 non-null   float64
7   alpha                                 8887 non-null   float64
8   Tide_water_level_cat_encoded          8887 non-null   int64
9   Tide_ebb_flood_cat_encoded            8887 non-null   int64
10  press_change_cat_encoded              8887 non-null   int64
11  Light                                 8887 non-null   object
12  tempF_Diff_encoded                    8887 non-null   int64
13  press                                 8887 non-null   float64
14  Moon_Phase_cat                        8887 non-null   object
dtypes: float64(8), int64(5), object(2)
memory usage: 1.0+ MB
```

In the next step, we do a rough feature selection by analyzing the information gain and chi-squared stats between features and targets.

Interpretation:

- Mutual Information tells you how much knowing a feature reduces uncertainty about the target.
- Chi-Square tells you whether there is a statistically significant relationship between a categorical feature and the target.

Practical Explanation Using Your Code

- If "Mutual Information" for Light is high, it means Light helps predict #sharks well.
- If "Chi-Square" for Light is significant, it means Light and #sharks are not independent, but it doesn't tell us how much it improves prediction. Chi-Square goes from 0 to infinity - the larger the value the more the dependence between the features.

Mutual information is better for feature selection in machine learning, while chi-square is useful for statistical hypothesis testing.

Example: Let's say we're analyzing Education Level vs. Income Group:

Chi-Square might only tell us if there's a relationship.

Mutual Information tells us how much knowing Education Level helps predict Income Group.

Feature Importance for model can be explained globally or locally:

Global Feature Importance: How important each feature is for the entire dataset.

Local Feature Importance: How important each feature is for a single prediction.

```

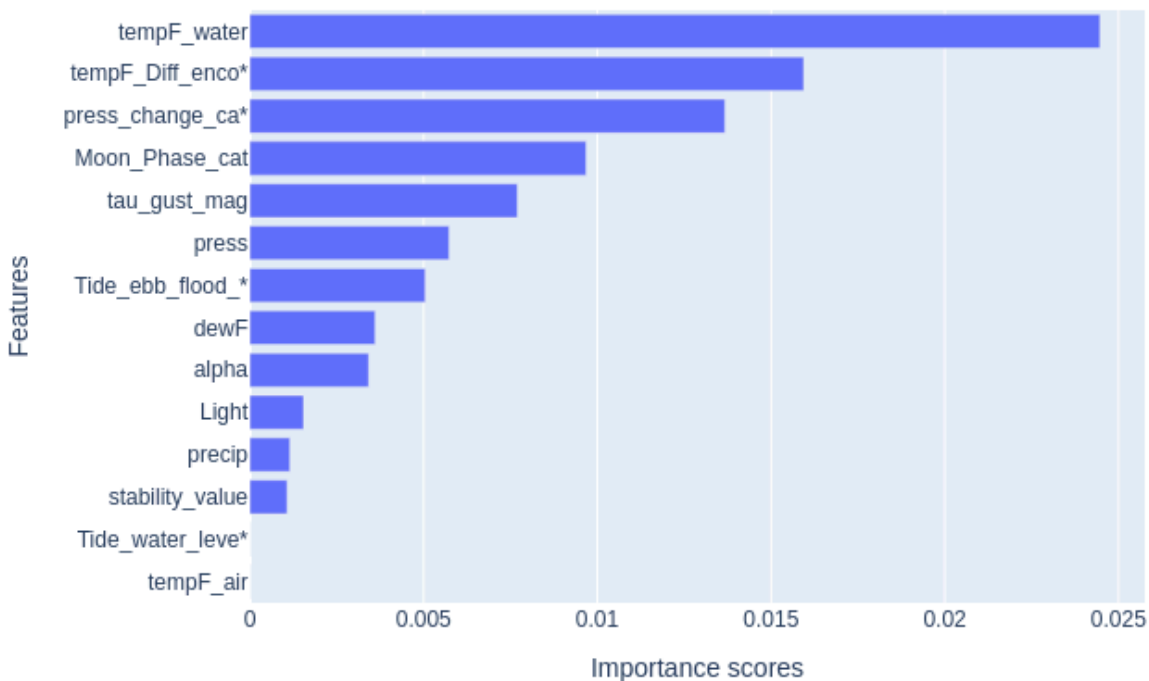
1 tabular_data = Tabular(
2     df2,
3     categorical_columns=['Light', 'Moon_Phase_cat'],
4     target_column='#sharks'
5 )
6 explainer = DataAnalyzer(
7     explainers=["mutual", "chi2"],
8     mode="classification",
9     data=tabular_data
10 )
11 data_explanations = explainer.explain_global()
12
13 print("Mutual information:")
14 data_explanations["mutual"].ipython_plot()
15 print("Chi square:")
16 data_explanations["chi2"].ipython_plot()

```



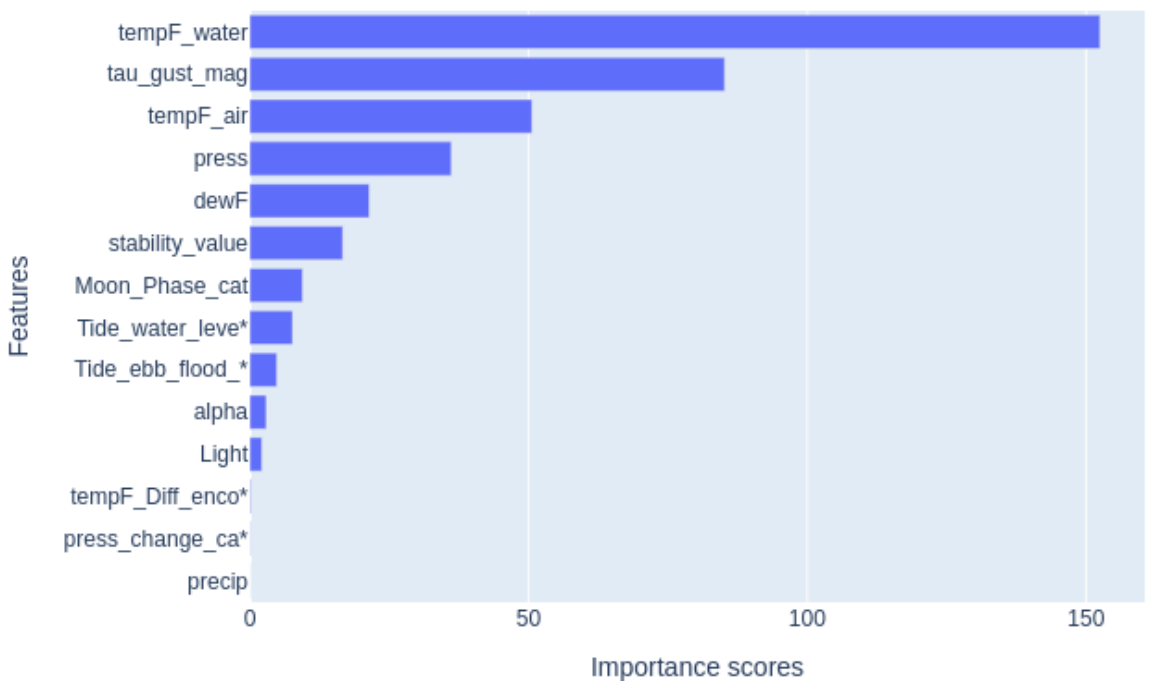

Mutual information:

Global Feature Importance



Chi square:

Global Feature Importance



The most important features showed above are "tempF_water", "tau_gust_mag", "tempF_air", "press", "dewF", "stability_value", "Moon_Phase_cat", "Tide_water_level", "Tide_ebb_flood_cat", and "alpha". Domain information will guide us to keep all but tempF_Diff_encoded" and "press_change_cat" in the dataframe for modeling.

```
1 # We drop this features because the feature has relatively low importance scores.
2 df3 = df2.drop(columns=["tempF_Diff_encoded"])
3 tabular_data = Tabular(
4     df3,
5     categorical_columns=['Light', 'Moon_Phase_cat'],
6     target_column='#sharks'
7 )
8 #print(tabular_data)
```

Train two models:

1. XGBoost classifier
2. Random Forest Classifier

Training XGBoost and Random Forest together can be useful because they have different strengths and weaknesses. Using both allows you to compare performance, interpret feature importance, and potentially combine their outputs for better predictions.

Factor	When to Use RF	When to Use XGBoost
Dataset Size	Small to medium	Large datasets
Missing Data	Handles missing values well	Requires imputation
Overfitting Risk	Lower	Higher, but tunable
Training Time	Faster on small datasets	Faster on large datasets
Interpretability	Easier to explain	Harder to interpret
Performance on Complex Data	Good, but may miss subtle patterns	Captures complex relationships better


```
1 import sklearn
2 import xgboost
3 from omnixai.preprocessing.tabular import TabularTransform
4
5 np.random.seed(12345)
6 # Train an XGBoost model
7 transformer = TabularTransform().fit(tabular_data)
8 x = transformer.transform(tabular_data)
9 train, test, train_labels, test_labels = sklearn.model_selection.train_test_split(
10     x[:, :-1], x[:, -1], train_size=0.80, stratify=x[:, -1]
11 )
12 print('Training data shape: {}'.format(train.shape))
13 print('Test data shape:      {}'.format(test.shape))
14
15 class_names = transformer.class_names
16 #
17 gbtree = xgboost.XGBClassifier(
18     objective="binary:logistic",
19     n_estimators=700,      # More trees for better learning
20     learning_rate=0.1,    # Lower learning rate to generalize better
21     max_depth=7,          # Slightly deeper trees
22     subsample=0.9,        # Use 85% of data per tree
23     colsample_bytree=1,   # Use 80% of features per tree
24     gamma=2,              # Avoid unnecessary splits
25     reg_lambda=3,         # Add regularization
26     reg_alpha=1           # Helps with feature selection
27 )
28 gbtree.fit(train, train_labels)
29 print('Test accuracy: {}'.format(
30     sklearn.metrics.accuracy_score(test_labels, gbtree.predict(test)))
31 # Convert the transformed data back to Tabular instances
32 train_data = transformer.invert(train)
33 test_data = transformer.invert(test)
```

➡ Training data shape: (7109, 23)
Test data shape: (1778, 23)
Test accuracy: 0.6827896512935883

```

1 np.random.seed(12345)
2 # Train an Random Forest model
3 transformer = TabularTransform().fit(tabular_data)
4 x = transformer.transform(tabular_data)
5 train, test, train_labels, test_labels = sklearn.model_selection.train_test_split(
6     x[:, :-1], x[:, -1], train_size=0.80, stratify=x[:, -1]
7 )
8 print('Training data shape: {}'.format(train.shape))
9 print('Test data shape:     {}'.format(test.shape))
10
11 class_names = transformer.class_names
12 # Train Random Forest on your training data
13 rf_model = RandomForestClassifier(
14     n_estimators=1000,      # Increase trees for better learning
15     max_depth=12,          # Limit tree depth to avoid overfitting
16     min_samples_split=4,   # Prevent too many small splits
17     min_samples_leaf=3,    # Avoid overfitting small leaf nodes
18     max_features=.8,       # Select a subset of features per tree
19     class_weight="balanced", # Handle class imbalance
20     random_state=42,
21     n_jobs=-1 # Use all CPU cores for faster training
22 )
23 # The line below has been modified to use the transformed data 'train' instead of 'train_data'
24 rf_model.fit(train, train_labels) # Use the transformed data for training
25 #
26 print('Test accuracy: {}'.format(
27     sklearn.metrics.accuracy_score(test_labels, rf_model.predict(test))))
28
29 # Convert the transformed data back to Tabular instances
30 train_data = transformer.invert(train)
31 test_data = transformer.invert(test)

```

 Training data shape: (7109, 23)
 Test data shape: (1778, 23)
 Test accuracy: 0.6912260967379078

We then create an `TabularExplainer` explainer to generate local and global explanations.

```

1 from omnixai.explainers.tabular import TabularExplainer
2
3 # Initialize a TabularExplainer
4 explainers = TabularExplainer(
5     explainers=["lime", "shap", "mace", "pdp", "ale"],
6     mode="classification",
7     data=train_data,
8     model=rf_model,
9     preprocess=lambda z: transformer.transform(z),
10    params={
11        "lime": {"kernel_width": 3},
12        "shap": {"nsamples": 100},
13        "mace": {"ignored_features": ["tempF_Diff_encoded"]}
14    }
15 )

```

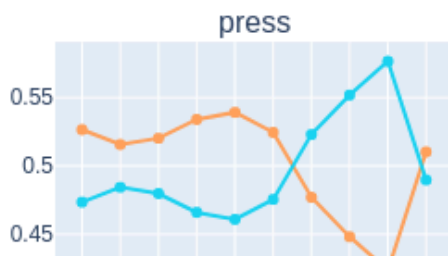
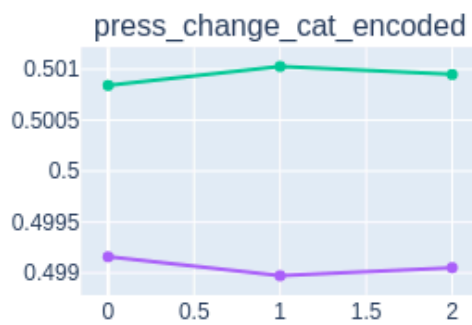
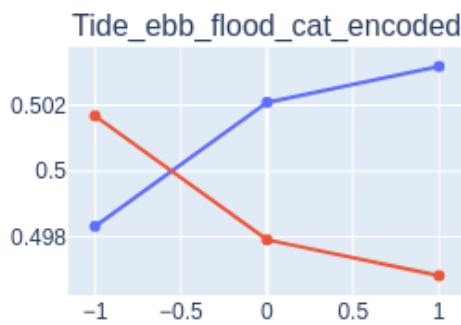
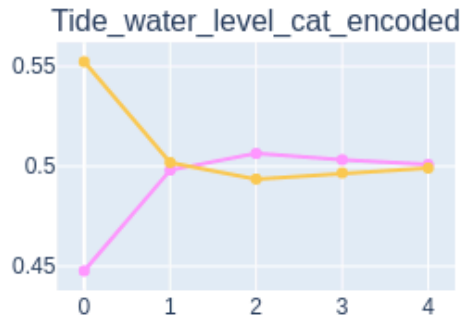
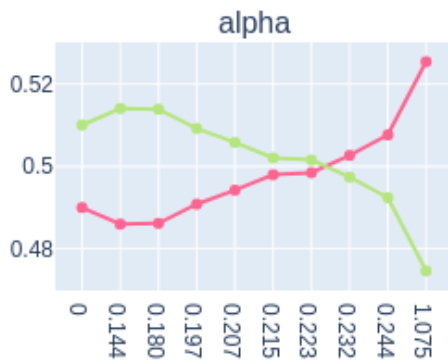
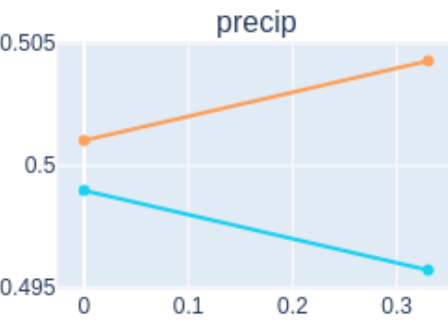
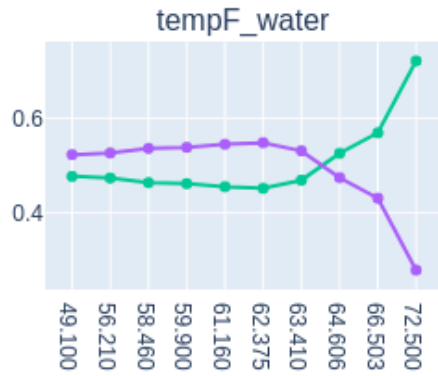
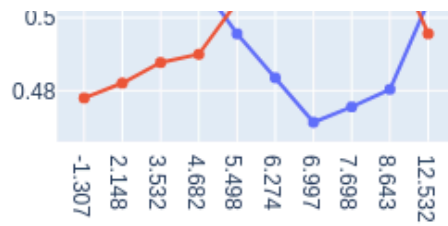
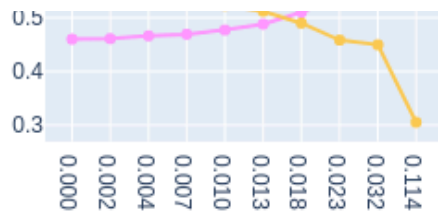
Interpretation of first chart: for Dusk, the model strongly favors Class = 1 (Red Line for Class 1) rather than Class 0 (Blue Line for Class 0)

You can plot just Class 1 line to more easily interpret the charts.

```

1 # Generate global explanations
2 global_explanations = explainers.explain_global()
3 global_explanations["pdp"].ipython_plot(class_names=class_names)
4 #The separation between the two lines shows that the model sees a clear difference

```





PDP compared to ALE

Method

PDP Example: How Does "Light" (Dawn, Day, Dusk, Night) Affect Predictions? PDP shows how changing "Light" affects the probability of Class 1. If "Night" has a **higher probability of Class 1**, it suggests that being at night increases the likelihood of a shark presence.

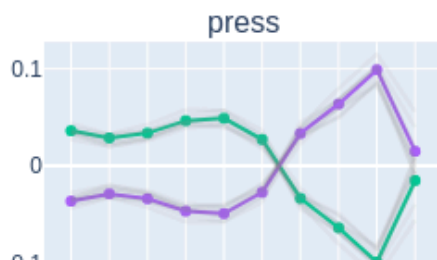
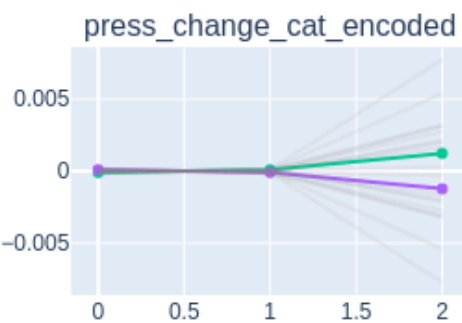
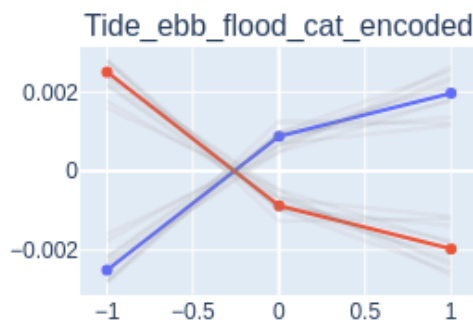
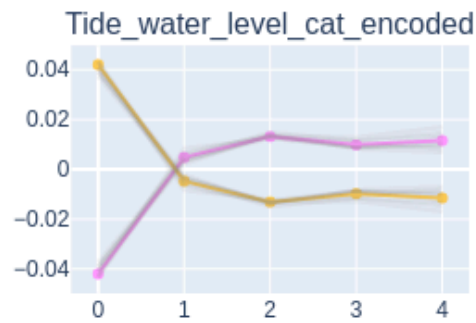
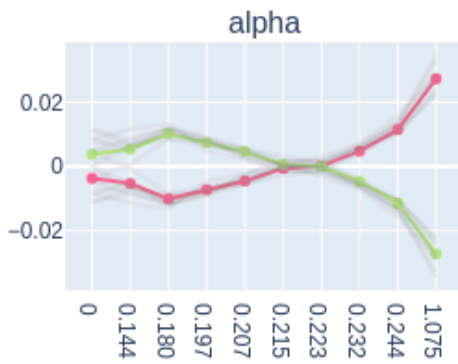
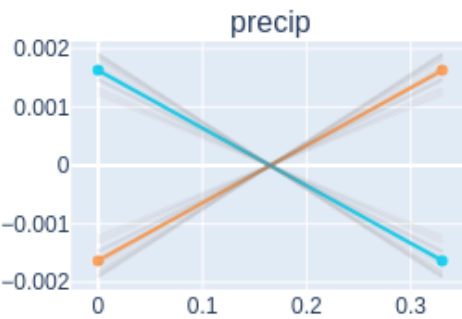
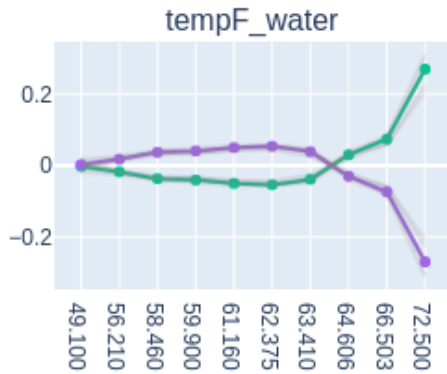
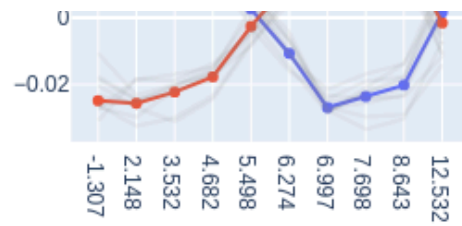
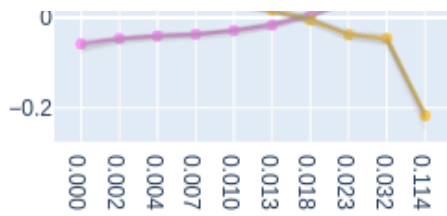
Method

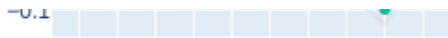
Problem: If "Light" is correlated with "Moon Phase" (e

ALE Example: How Does "Light" Affect Predictions?

ALE looks at **small changes** in "Light" and their direct
If "Light" is correlated with "Moon Phase", ALE adjusts
Instead of showing a global trend, ALE focuses on **ho**

```
1 global_explanations["ale"].ipython_plot(class_names=class_names)
```



For some specific test instances, we can generate local explanations to analyze the predictions.

```
1 num_records = len(test_data)
2 print("Number of records in test_data:", num_records)
```