

Anomaly detection in unsupervised learning involves identifying rare items, events, or observations that do not conform to the general pattern of the data. These anomalies can indicate critical incidents such as technical glitches, fraud, or diseases. Since the data is unlabeled, unsupervised anomaly detection methods rely on inherent properties of the data to detect outliers. Here are some common anomaly detection algorithms used in unsupervised learning:

Isolation Forest

is an unsupervised learning algorithm primarily used for anomaly detection. It works on the principle that anomalies are "few and different" from the majority of the data. This algorithm isolates observations by randomly selecting a feature and then randomly selecting a split value between the maximum and minimum values of the selected feature. The idea is that anomalies are more likely to be isolated closer to the root of the tree.

Purpose

- **Anomaly Detection:** Identifying outliers or anomalies in the data.
- **Fraud Detection:** Detecting fraudulent transactions.
- **Network Security:** Identifying unusual patterns in network traffic.
- **Fault Detection:** Detecting faults in industrial systems.

Suitable Data Types

- **Numerical Data:** Continuous or discrete numerical values.
- **High-dimensional Data:** Effective even with a large number of features.
- **Mixed Data:** Can handle datasets with both categorical and numerical features (though typically requires preprocessing for categorical features).

Examples

1. **Credit Card Fraud Detection:** Identifying fraudulent transactions from a dataset of credit card transactions.
2. **Network Intrusion Detection:** Detecting unusual network traffic patterns indicating potential security breaches.
3. **Healthcare:** Identifying rare diseases or abnormal patient data.
4. **Manufacturing:** Detecting defective products or equipment failures.

Implementation and Visualization

Let's implement the Isolation Forest algorithm using a real-life dataset and visualize the data before and after applying the model.

We'll use the "Credit Card Fraud Detection" dataset from Kaggle for this demonstration. This dataset contains transactions made by credit cards in September 2013 by European cardholders. The dataset presents transactions that occurred in two days, where we have 492 frauds out of 284,807 transactions.

First, let's load the dataset and visualize the data before applying the Isolation Forest model.

```
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.ensemble import IsolationForest

# Load the dataset
data = pd.read_csv('creditcard.csv')

# Visualize the distribution of the 'Amount' feature
plt.figure(figsize=(10, 6))
plt.hist(data['Amount'], bins=50, color='blue', edgecolor='black')
plt.title('Transaction Amount Distribution')
plt.xlabel('Amount')
plt.ylabel('Frequency')
plt.show()

# Visualize the data before applying Isolation Forest
plt.figure(figsize=(10, 6))
plt.scatter(data['Time'], data['Amount'], c='blue', alpha=0.5)
plt.title('Transaction Amount over Time')
plt.xlabel('Time')
plt.ylabel('Amount')
plt.show()

Next, we'll apply the Isolation Forest algorithm and visualize the results.

# Define the model
model = IsolationForest(contamination=0.001)

# Fit the model
data['anomaly'] = model.fit_predict(data[['Time', 'Amount']])
```

```

# Visualize the anomalies

anomalies = data[data['anomaly'] == -1]
normal = data[data['anomaly'] == 1]

plt.figure(figsize=(10, 6))

plt.scatter(normal['Time'], normal['Amount'], c='blue', alpha=0.5, label='Normal')

plt.scatter(anomalies['Time'], anomalies['Amount'], c='red', alpha=0.5, label='Anomaly')

plt.title('Transaction Amount over Time (After Isolation Forest)')

plt.xlabel('Time')

plt.ylabel('Amount')

plt.legend()

plt.show()

```

Code Explanation

- **import pandas as pd:** Imports the pandas library, which is used for data manipulation and analysis.
 - **import matplotlib.pyplot as plt:** Imports the pyplot module from matplotlib, which is used for data visualization.
 - **from sklearn.ensemble import IsolationForest:** Imports the IsolationForest class from the sklearn.ensemble module, which is used to create an Isolation Forest model for anomaly detection.
- data = pd.read_csv('creditcard.csv'):** Reads the CSV file containing the dataset into a pandas DataFrame called data.
- **plt.figure(figsize=(10, 6)):** Creates a new figure with a specified size (10 inches by 6 inches).
 - **plt.hist(data['Amount'], bins=50, color='blue', edgecolor='black'):** Plots a histogram of the 'Amount' feature in the dataset with 50 bins, blue bars, and black edges.
 - **plt.title('Transaction Amount Distribution'):** Sets the title of the plot to 'Transaction Amount Distribution'.
 - **plt.xlabel('Amount'):** Sets the label for the x-axis to 'Amount'.
 - **plt.ylabel('Frequency'):** Sets the label for the y-axis to 'Frequency'.
 - **plt.show():** Displays the plot.
 - **plt.figure(figsize=(10, 6)):** Creates a new figure with a specified size (10 inches by 6 inches).
 - **plt.scatter(data['Time'], data['Amount'], c='blue', alpha=0.5):** Creates a scatter plot of the 'Time' feature on the x-axis and the 'Amount' feature on the y-axis. Points are colored blue and have an alpha value of 0.5 (transparency).
 - **plt.title('Transaction Amount over Time'):** Sets the title of the plot to 'Transaction Amount over Time'.

- ❑ **plt.xlabel('Time')**: Sets the label for the x-axis to 'Time'.
- ❑ **plt.ylabel('Amount')**: Sets the label for the y-axis to 'Amount'.
- ❑ **plt.show()**: Displays the plot.

model = IsolationForest(contamination=0.001): Creates an instance of the IsolationForest model with a contamination parameter of 0.001, indicating that we expect 0.1% of the data to be anomalies.

data['anomaly'] = model.fit_predict(data[['Time', 'Amount']]): Fits the IsolationForest model on the 'Time' and 'Amount' features of the data and adds a new column 'anomaly' to the DataFrame, where the model predicts -1 for anomalies and 1 for normal points.

- ❑ **anomalies = data[data['anomaly'] == -1]**: Creates a DataFrame anomalies containing only the rows where the 'anomaly' column is -1 (predicted anomalies).
- ❑ **normal = data[data['anomaly'] == 1]**: Creates a DataFrame normal containing only the rows where the 'anomaly' column is 1 (predicted normal points).
- ❑ **plt.figure(figsize=(10, 6))**: Creates a new figure with a specified size (10 inches by 6 inches).
- ❑ **plt.scatter(normal['Time'], normal['Amount'], c='blue', alpha=0.5, label='Normal')**: Creates a scatter plot of the normal data points with 'Time' on the x-axis and 'Amount' on the y-axis. Points are colored blue, have an alpha value of 0.5 (transparency), and are labeled 'Normal'.
- ❑ **plt.scatter(anomalies['Time'], anomalies['Amount'], c='red', alpha=0.5, label='Anomaly')**: Creates a scatter plot of the anomaly data points with 'Time' on the x-axis and 'Amount' on the y-axis. Points are colored red, have an alpha value of 0.5 (transparency), and are labeled 'Anomaly'.
- ❑ **plt.title('Transaction Amount over Time (After Isolation Forest)')**: Sets the title of the plot to 'Transaction Amount over Time (After Isolation Forest)'.
- ❑ **plt.xlabel('Time')**: Sets the label for the x-axis to 'Time'.
- ❑ **plt.ylabel('Amount')**: Sets the label for the y-axis to 'Amount'.
- ❑ **plt.legend()**: Adds a legend to the plot to differentiate between normal and anomaly points.
- ❑ **plt.show()**: Displays the plot.