

```

14290           NaN           NaN
14291           NaN           NaN
14292           NaN           NaN
14293           NaN           NaN
14294           NaN           NaN
14295           NaN           NaN

   Gyro_Magnitude_Combined_Outliers  label_Combined_Outliers
14286                  NaN           NaN
14287                  NaN           NaN
14288                  NaN           NaN
14289                  NaN           NaN
14290                  NaN           NaN
14291                  NaN           NaN
14292                  NaN           NaN
14293                  NaN           NaN
14294                  NaN           NaN
14295                  NaN           NaN

```

[10 rows x 53 columns]

Explanation of Changes:

`synthetic_anomalies`: A list is used to collect synthetic anomaly rows.

`pd.DataFrame(synthetic_anomalies)`: Converts the list of anomalies into a DataFrame.

`pd.concat()`: Combines the original DataFrame (`df`) with the new DataFrame of synthetic anomalies.

6a: Isolation Forest Model Setup

- i. Experiment with contamination rates: Test different contamination rates (e.g., 0.01, 0.05, 0.1) and document the number of anomalies detected.
- ii. Test different number of estimators: Vary the number of base estimators in the ensemble (e.g., 50, 100, 200) and analyze their impact.
- iii. Visualize anomaly scores: Plot anomaly scores for different hyperparameters and interpret the results.
- iv. Analyze model performance for each hyperparameter set: Compare how different hyperparameter combinations affect the isolation forest's detection accuracy.

```
[48]: from sklearn.ensemble import IsolationForest

# Contamination rates to test
contamination_rates = [0.01, 0.05, 0.1]

# Dictionary to store results for each contamination rate
anomalies_detected = {}

for rate in contamination_rates:
```

```

# Initialize and fit Isolation Forest with a specific contamination rate
iso_forest = IsolationForest(contamination=rate, random_state=42)
iso_forest.fit(X)

# Predict anomalies (-1 for anomaly, 1 for normal)
y_pred = iso_forest.predict(X)
anomalies_count = np.sum(y_pred == -1) # Count anomalies detected

# Store results
anomalies_detected[rate] = anomalies_count
print(f"Contamination Rate: {rate}, Anomalies Detected: {anomalies_count}")

# Display results summary
print("\nSummary of anomalies detected for different contamination rates:")
for rate, count in anomalies_detected.items():
    print(f"Contamination Rate: {rate} - Anomalies Detected: {count}")

```

Contamination Rate: 0.01, Anomalies Detected: 143
Contamination Rate: 0.05, Anomalies Detected: 713
Contamination Rate: 0.1, Anomalies Detected: 1425

Summary of anomalies detected for different contamination rates:
Contamination Rate: 0.01 - Anomalies Detected: 143
Contamination Rate: 0.05 - Anomalies Detected: 713
Contamination Rate: 0.1 - Anomalies Detected: 1425

```

[49]: # Estimator values to test
estimator_counts = [50, 100, 200]

# Dictionary to store results for each number of estimators
anomalies_by_estimators = {}

for n_estimators in estimator_counts:
    # Initialize and fit Isolation Forest with a specific number of estimators
    iso_forest = IsolationForest(n_estimators=n_estimators, contamination=0.05,random_state=42)
    iso_forest.fit(X)

    # Predict anomalies
    y_pred = iso_forest.predict(X)
    anomalies_count = np.sum(y_pred == -1) # Count anomalies detected

    # Store results
    anomalies_by_estimators[n_estimators] = anomalies_count
    print(f"Estimators: {n_estimators}, Anomalies Detected: {anomalies_count}")

# Display results summary

```

```

print("\nSummary of anomalies detected for different numbers of estimators:")
for n_estimators, count in anomalies_by_estimators.items():
    print(f"Estimators: {n_estimators} - Anomalies Detected: {count}")

```

Estimators: 50, Anomalies Detected: 713
 Estimators: 100, Anomalies Detected: 713
 Estimators: 200, Anomalies Detected: 713

Summary of anomalies detected for different numbers of estimators:
 Estimators: 50 - Anomalies Detected: 713
 Estimators: 100 - Anomalies Detected: 713
 Estimators: 200 - Anomalies Detected: 713

```

[50]: import matplotlib.pyplot as plt

# Function to plot anomaly scores
def plot_anomaly_scores(contamination, n_estimators):
    # Initialize Isolation Forest with specified hyperparameters
    iso_forest = IsolationForest(contamination=contamination,
                                  n_estimators=n_estimators, random_state=42)
    iso_forest.fit(X)

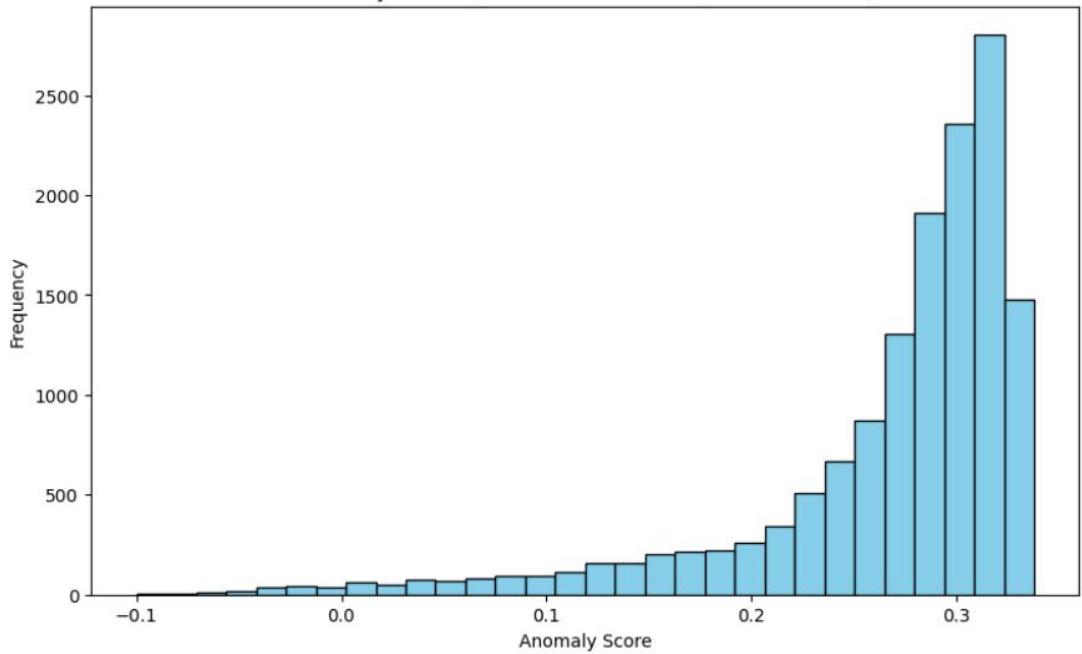
    # Retrieve anomaly scores (lower scores indicate higher likelihood of being
    # an anomaly)
    anomaly_scores = iso_forest.decision_function(X)  # Higher values mean less
    # anomalous
    plt.figure(figsize=(10, 6))

    # Plot anomaly scores
    plt.hist(anomaly_scores, bins=30, color='skyblue', edgecolor='black')
    plt.title(f'Anomaly Scores (Contamination: {contamination}, Estimators: {n_estimators})')
    plt.xlabel('Anomaly Score')
    plt.ylabel('Frequency')
    plt.show()

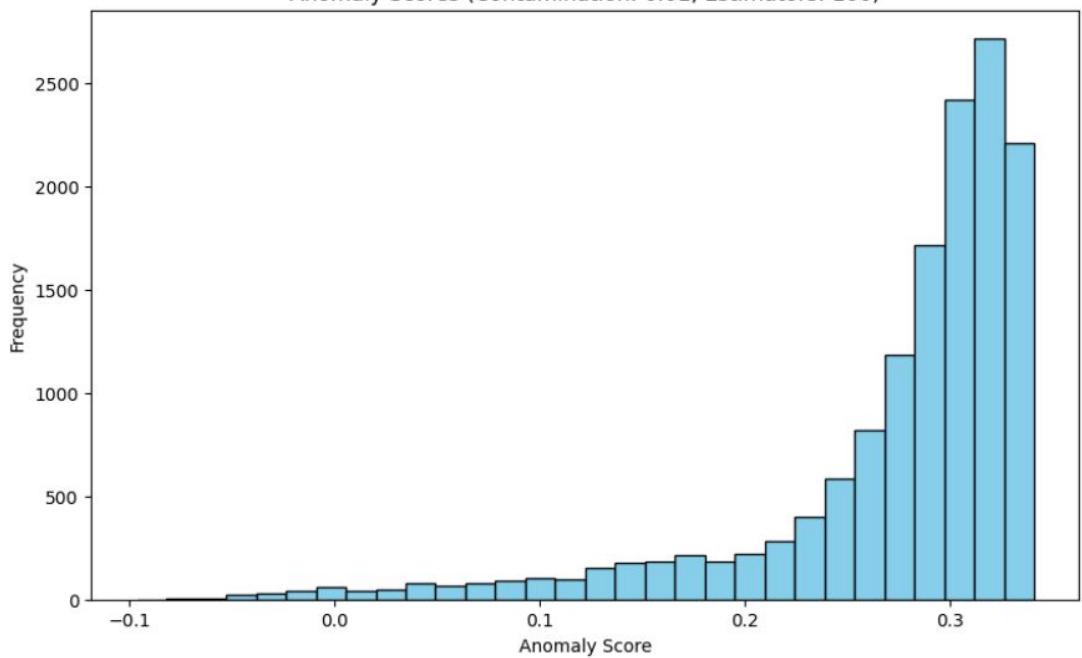
# Plot anomaly scores for different parameter combinations
for rate in contamination_rates:
    for n_estimators in estimator_counts:
        plot_anomaly_scores(rate, n_estimators)

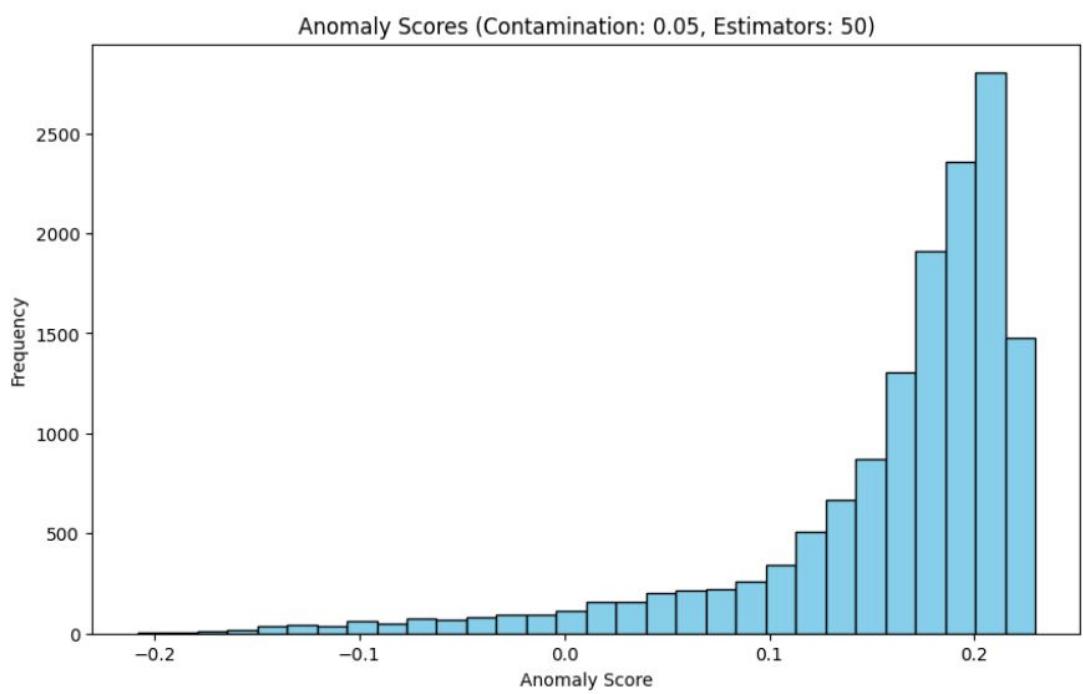
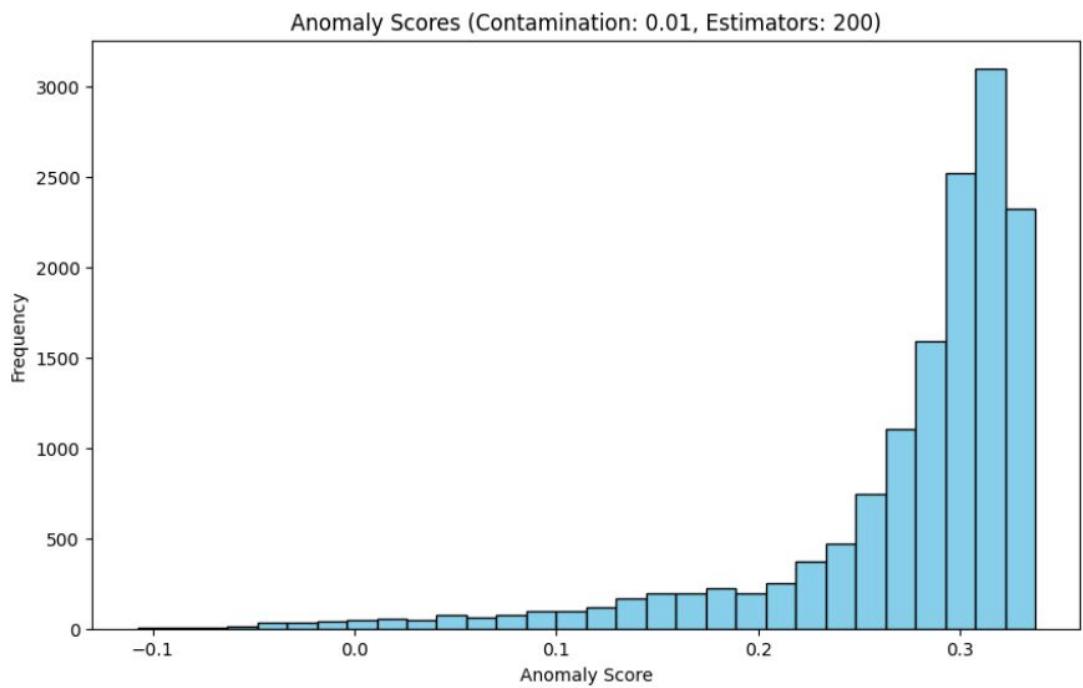
```

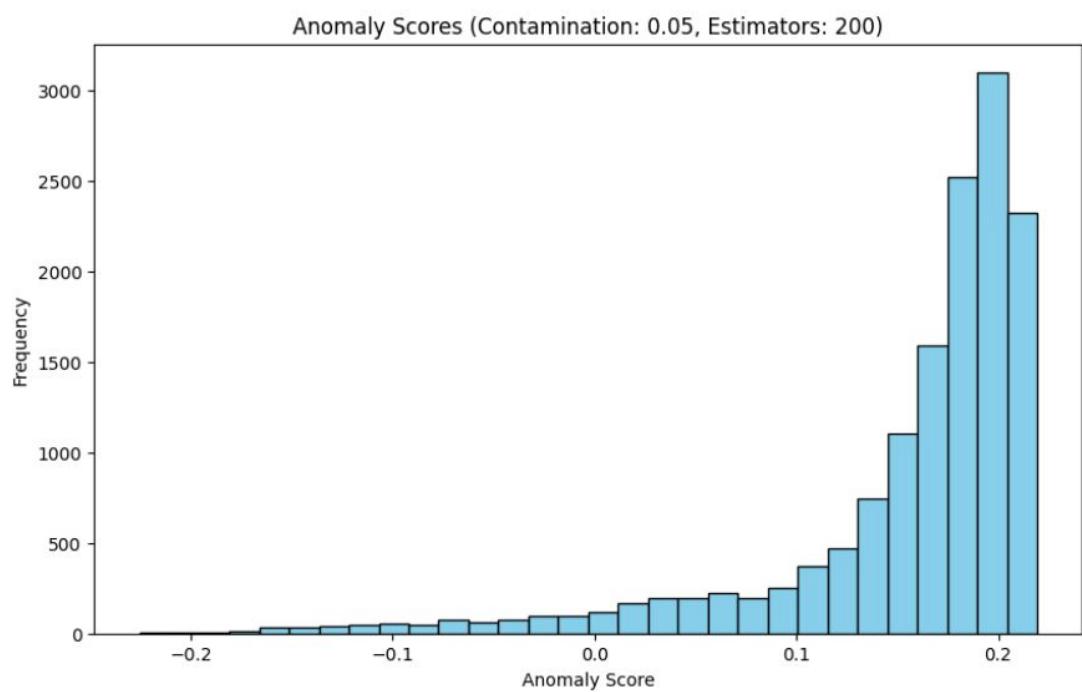
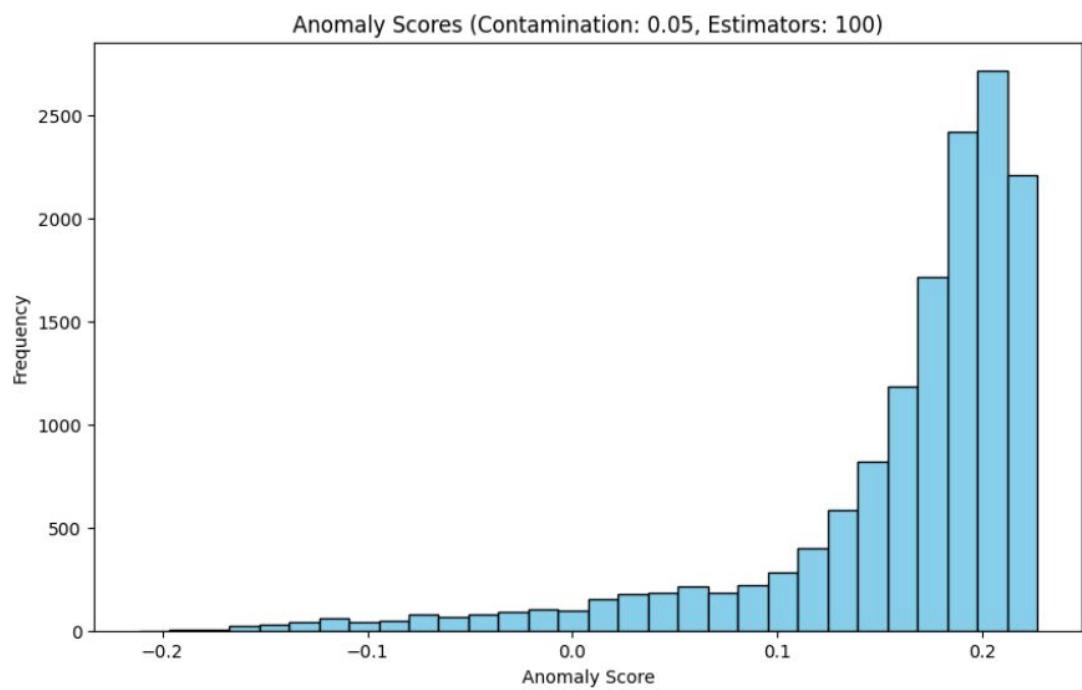
Anomaly Scores (Contamination: 0.01, Estimators: 50)

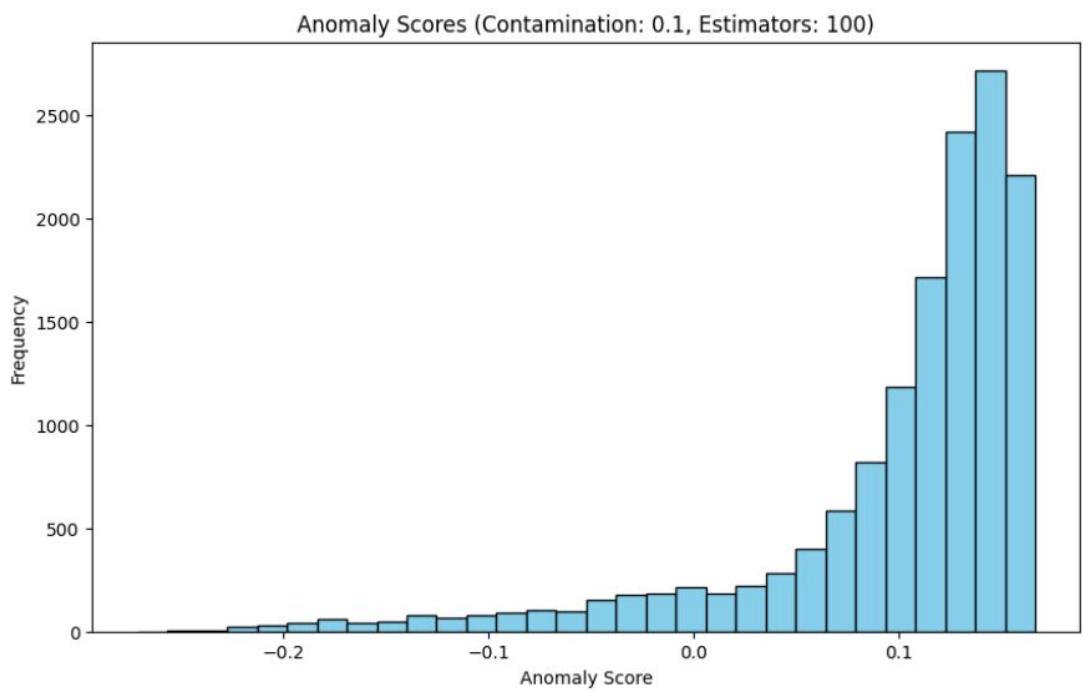
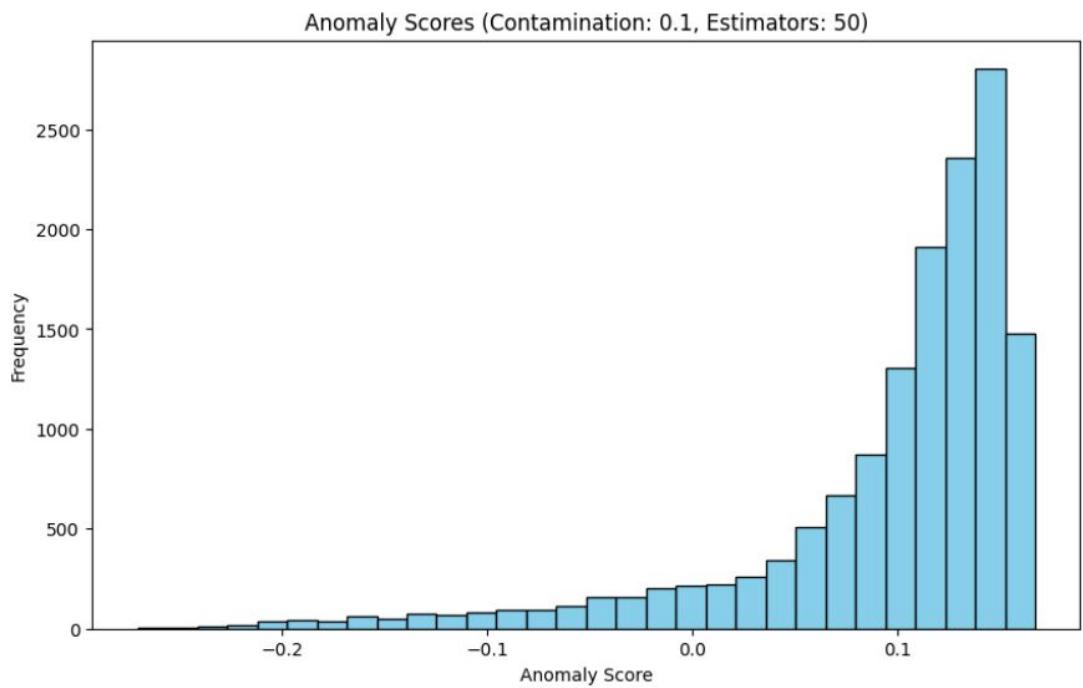


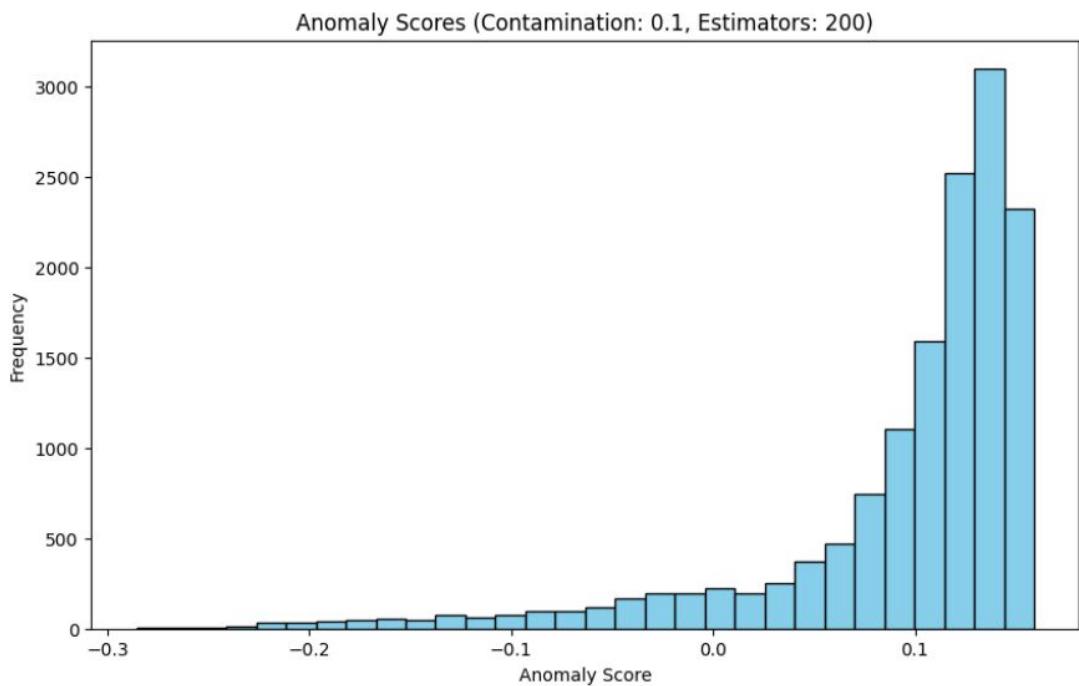
Anomaly Scores (Contamination: 0.01, Estimators: 100)











```
[51]: from sklearn.metrics import f1_score, precision_score, recall_score,
       accuracy_score, roc_auc_score

# Dictionary to store performance results for each hyperparameter combination
performance_results = {}

# Iterate over each combination of contamination rate and number of estimators
for rate in contamination_rates:
    for n_estimators in estimator_counts:
        # Initialize Isolation Forest with specific hyperparameters
        iso_forest = IsolationForest(contamination=rate,
                                      n_estimators=n_estimators, random_state=42)
        iso_forest.fit(X)

        # Predict and convert predictions to binary format (1 for anomaly, 0 for normal)
        y_pred = iso_forest.predict(X)
        y_pred = np.where(y_pred == -1, 1, 0) # Convert -1 to 1 (anomaly)

        # Calculate performance metrics
        f1 = f1_score(df['label'], y_pred)
        precision = precision_score(df['label'], y_pred)
        recall = recall_score(df['label'], y_pred)
        accuracy = accuracy_score(df['label'], y_pred)
```

```

roc_auc = roc_auc_score(df['label'], y_pred)

# Store performance metrics in dictionary
performance_results[(rate, n_estimators)] = {
    'F1-Score': f1,
    'Precision': precision,
    'Recall': recall,
    'Accuracy': accuracy,
    'ROC-AUC': roc_auc
}

# Display metrics for the current hyperparameter combination
print(f"Contamination Rate: {rate}, Estimators: {n_estimators}")
print(f"  F1-Score: {f1}")
print(f"  Precision: {precision}")
print(f"  Recall: {recall}")
print(f"  Accuracy: {accuracy}")
print(f"  ROC-AUC: {roc_auc}\n")

# Summary of all performance metrics for each combination
print("\nSummary of performance metrics for each hyperparameter combination:")
for params, metrics in performance_results.items():
    print(f"Contamination Rate: {params[0]}, Estimators: {params[1]}")
    for metric_name, metric_value in metrics.items():
        print(f"  {metric_name}: {metric_value}")
    print()

```

Contamination Rate: 0.01, Estimators: 50
 F1-Score: 0.03284799068142109
 Precision: 0.986013986013986
 Recall: 0.016702203269367447
 Accuracy: 0.4171697318545557
 ROC-AUC: 0.5081788066656968

Contamination Rate: 0.01, Estimators: 100
 F1-Score: 0.03284799068142109
 Precision: 0.986013986013986
 Recall: 0.016702203269367447
 Accuracy: 0.4171697318545557
 ROC-AUC: 0.5081788066656968

Contamination Rate: 0.01, Estimators: 200
 F1-Score: 0.03284799068142109
 Precision: 0.986013986013986
 Recall: 0.016702203269367447
 Accuracy: 0.4171697318545557
 ROC-AUC: 0.5081788066656968

Contamination Rate: 0.05, Estimators: 50

F1-Score: 0.14877116329874385
Precision: 0.9551192145862553
Recall: 0.0806680881307747
Accuracy: 0.4529692545275867
ROC-AUC: 0.5375773245615969

Contamination Rate: 0.05, Estimators: 100

F1-Score: 0.14898962315674494
Precision: 0.9565217391304348
Recall: 0.08078654347311064
Accuracy: 0.453109644812579
ROC-AUC: 0.5377226997172583

Contamination Rate: 0.05, Estimators: 200

F1-Score: 0.15008192244675042
Precision: 0.9635343618513323
Recall: 0.08137882018479034
Accuracy: 0.45381159623754036
ROC-AUC: 0.5384495754955654

Contamination Rate: 0.1, Estimators: 50

F1-Score: 0.2699908786865309
Precision: 0.9347368421052632
Recall: 0.15778251599147122
Accuracy: 0.49438438860030887
ROC-AUC: 0.5708795419378445

Contamination Rate: 0.1, Estimators: 100

F1-Score: 0.26938279112192154
Precision: 0.9326315789473684
Recall: 0.1574271499644634
Accuracy: 0.49396321774533203
ROC-AUC: 0.5704434164708602

Contamination Rate: 0.1, Estimators: 200

F1-Score: 0.2699908786865309
Precision: 0.9347368421052632
Recall: 0.15778251599147122
Accuracy: 0.49438438860030887
ROC-AUC: 0.5708795419378445

Summary of performance metrics for each hyperparameter combination:

Contamination Rate: 0.01, Estimators: 50

F1-Score: 0.03284799068142109
Precision: 0.986013986013986

Recall: 0.016702203269367447
Accuracy: 0.4171697318545557
ROC-AUC: 0.5081788066656968

Contamination Rate: 0.01, Estimators: 100
F1-Score: 0.03284799068142109
Precision: 0.986013986013986
Recall: 0.016702203269367447
Accuracy: 0.4171697318545557
ROC-AUC: 0.5081788066656968

Contamination Rate: 0.01, Estimators: 200
F1-Score: 0.03284799068142109
Precision: 0.986013986013986
Recall: 0.016702203269367447
Accuracy: 0.4171697318545557
ROC-AUC: 0.5081788066656968

Contamination Rate: 0.05, Estimators: 50
F1-Score: 0.14877116329874385
Precision: 0.9551192145862553
Recall: 0.0806680881307747
Accuracy: 0.4529692545275867
ROC-AUC: 0.5375773245615969

Contamination Rate: 0.05, Estimators: 100
F1-Score: 0.14898962315674494
Precision: 0.9565217391304348
Recall: 0.08078654347311064
Accuracy: 0.453109644812579
ROC-AUC: 0.5377226997172583

Contamination Rate: 0.05, Estimators: 200
F1-Score: 0.15008192244675042
Precision: 0.9635343618513323
Recall: 0.08137882018479034
Accuracy: 0.45381159623754036
ROC-AUC: 0.5384495754955654

Contamination Rate: 0.1, Estimators: 50
F1-Score: 0.2699908786865309
Precision: 0.9347368421052632
Recall: 0.15778251599147122
Accuracy: 0.49438438860030887
ROC-AUC: 0.5708795419378445

Contamination Rate: 0.1, Estimators: 100
F1-Score: 0.26938279112192154

```

Precision: 0.9326315789473684
Recall: 0.1574271499644634
Accuracy: 0.49396321774533203
ROC-AUC: 0.5704434164708602

Contamination Rate: 0.1, Estimators: 200
F1-Score: 0.2699908786865309
Precision: 0.9347368421052632
Recall: 0.15778251599147122
Accuracy: 0.49438438860030887
ROC-AUC: 0.5708795419378445

```

6b: Performance Analysis for Isolation Forest

- i. Calculate precision, recall, and F1-score: After running the model, calculate these metrics using known anomalies in the data.
- ii. Plot precision-recall curves: Visualize the trade-off between precision and recall to help find the optimal model configuration.
- iii. Identify false positives and false negatives: Investigate cases where normal data is flagged as anomalous or anomalies are missed.
- iv. Compare results with previous methods: Compare Isolation Forest performance with earlier methods (IQR, Z-score) to determine which is more accurate.

```
[52]: from sklearn.metrics import precision_score, recall_score, f1_score

# Calculate precision, recall, and F1-score for Isolation Forest results
y_true = df['label'] # Known anomalies (true labels)
y_pred = df['IsoForest_Anomaly'] # Isolation Forest predictions

precision = precision_score(y_true, y_pred)
recall = recall_score(y_true, y_pred)
f1 = f1_score(y_true, y_pred)

print("Isolation Forest Performance Metrics:")
print(f" Precision: {precision}")
print(f" Recall: {recall}")
print(f" F1-Score: {f1}")
```

Isolation Forest Performance Metrics:

```

Precision: 0.9313358302122348
Recall: 0.1767353707652215
F1-Score: 0.29709279171644765

```

```
[53]: from sklearn.metrics import precision_recall_curve
import matplotlib.pyplot as plt
```

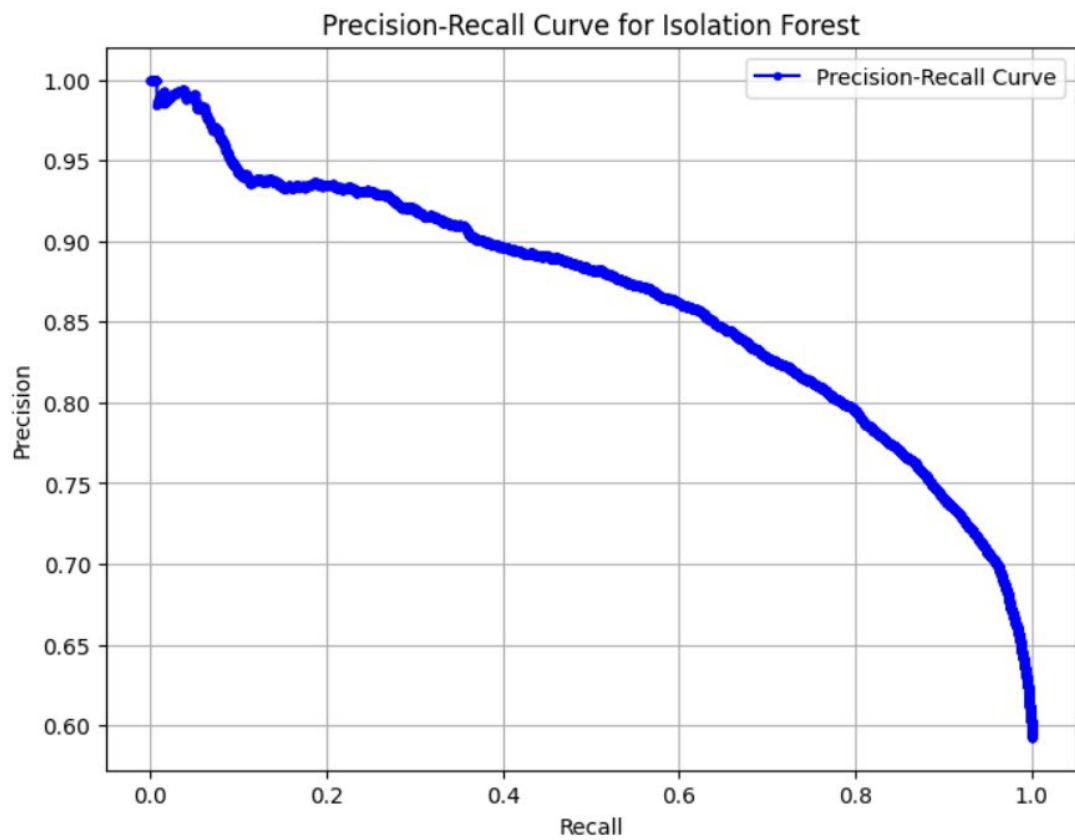
```

# Calculate anomaly scores for each instance in the dataset
anomaly_scores = -iso_forest.decision_function(X) # Higher scores indicate
# higher anomaly probability

# Calculate precision-recall curve
precision, recall, thresholds = precision_recall_curve(y_true, anomaly_scores)

# Plot the precision-recall curve
plt.figure(figsize=(8, 6))
plt.plot(recall, precision, marker='.', color='blue', label="Precision-Recall
#Curve")
plt.title("Precision-Recall Curve for Isolation Forest")
plt.xlabel("Recall")
plt.ylabel("Precision")
plt.legend()
plt.grid()
plt.show()

```



```
[54]: # Identify false positives and false negatives
false_positives = df[(y_pred == 1) & (y_true == 0)]
false_negatives = df[(y_pred == 0) & (y_true == 1)]

# Display counts and sample data
print(f"False Positives (Normal data flagged as anomalies): {len(false_positives)}")
print(false_positives.head())

print(f"\nFalse Negatives (Anomalies missed by model): {len(false_negatives)}")
print(false_negatives.head())
```

False Positives (Normal data flagged as anomalies): 110

	Longitude	Latitude	Speed	Distance	Acc X	Acc Y	Acc Z
943	73.821229	18.504085	0.000	311.2332	-1.126007	1.740718	0.338949
3853	73.828233	18.507553	1.554	881.1768	-1.283974	1.615290	-4.134312
5928	73.829893	18.516211	1.634	1785.8365	1.869302	-2.121189	2.896130
6024	73.830540	18.515993	1.566	1794.6466	-1.980159	-0.629080	3.044819
6025	73.830540	18.515993	1.566	1794.6466	-2.892073	-1.075845	3.499141

	Heading	gyro_x	gyro_y	hour	minute	second	mahalanobis
943	72.0	0.377250	0.078517	...	18	46	47
3853	117.0	-0.240335	-0.128567	...	18	51	38
5928	184.0	0.092587	0.102951	...	18	55	6
6024	160.0	-0.270267	0.181753	...	18	55	16
6025	160.0	-0.253774	0.142657	...	18	55	16

	mahalanobis_outlier	IsoForest_Anomaly	Total_Acc_Combined_Outliers
943	False	1	False
3853	False	1	False
5928	False	1	False
6024	False	1	False
6025	False	1	False

	Acc_Magnitude_Combined_Outliers	Gyro_Magnitude_Combined_Outliers
943	False	False
3853	False	False
5928	False	False
6024	False	False
6025	False	False

	label_Combined_Outliers
943	False
3853	False
5928	False
6024	False
6025	False

[5 rows x 53 columns]

```
False Negatives (Anomalies missed by model): 6950
   Longitude  Latitude  Speed  Distance    Acc X    Acc Y    Acc Z \
160  73.822520  18.501616  0.732  14.906174 -2.049965 -0.575076 -1.907734
161  73.822520  18.501616  0.732  14.906174 -0.193694  0.668475 -1.800664
175  73.822482  18.501607  0.820  18.999422 -1.346871  1.717449 -1.352394
176  73.822482  18.501607  0.820  18.999422 -0.679483  1.536472 -0.936674
177  73.822482  18.501607  0.820  18.999422 -0.664038  1.244627 -0.975004

   Heading  gyro_x  gyro_y ... hour minute second mahalanobis \
160  352.0 -0.262326 -0.146893 ... 18     45      28     1.427326
161  353.0  0.104804  0.068132 ... 18     45      28     0.829307
175  354.0 -0.028975  0.093177 ... 18     45      30     1.677286
176  353.0  0.010731  0.066299 ... 18     45      30     0.826869
177  354.0  0.115189  0.066910 ... 18     45      30     0.726693

   mahalanobis_outlier IsoForest_Anomaly Total_Acc_Combined_Outliers \
160           False                  0                   False
161           False                  0                   False
175           False                  0                   False
176           False                  0                   False
177           False                  0                   False

   Acc_Magnitude_Combined_Outliers Gyro_Magnitude_Combined_Outliers \
160                      False                      False
161                      False                      False
175                      False                      False
176                      False                      False
177                      False                      False

   label_Combined_Outliers
160           False
161           False
175           False
176           False
177           False
```

[5 rows x 53 columns]

```
[55]: # Comparison function for Isolation Forest, IQR, and Z-score methods
def compare_methods(df):
    # Isolation Forest metrics
    iso_precision = precision_score(y_true, y_pred)
    iso_recall = recall_score(y_true, y_pred)
    iso_f1 = f1_score(y_true, y_pred)
```

```

# IQR metrics
iqr_pred = df['label_Total_Acc_IQR_Outliers']
iqr_precision = precision_score(y_true, iqr_pred)
iqr_recall = recall_score(y_true, iqr_pred)
iqr_f1 = f1_score(y_true, iqr_pred)

# Z-score metrics
zscore_pred = df['label_Total_Acc_Z_Outliers']
zscore_precision = precision_score(y_true, zscore_pred)
zscore_recall = recall_score(y_true, zscore_pred)
zscore_f1 = f1_score(y_true, zscore_pred)

print("Comparison of Outlier Detection Methods:")
print(f"Isolation Forest - Precision: {iso_precision}, Recall: {iso_recall}, F1-Score: {iso_f1}")
print(f"IQR - Precision: {iqr_precision}, Recall: {iqr_recall}, F1-Score: {iqr_f1}")
print(f"Z-score - Precision: {zscore_precision}, Recall: {zscore_recall}, F1-Score: {zscore_f1}")

# Call the comparison function
compare_methods(df)

```

Comparison of Outlier Detection Methods:

Isolation Forest - Precision: 0.9313358302122348, Recall: 0.1767353707652215, F1-Score: 0.29709279171644765

IQR - Precision: 0.9677790563866513, Recall: 0.099620942904525, F1-Score: 0.18064654709483408

Z-score - Precision: 0.986013986013986, Recall: 0.033404406538734895, F1-Score: 0.06461961503208066

6c: Hyperparameter Tuning for Isolation Forest

- Test different contamination rates: Start with a small range of contamination rates (e.g., 0.01 to 0.1) and record their impact on outlier detection.
- Vary the max_samples parameter: Experiment with different max_samples settings (e.g., 0.5, 0.75, 1.0) to test model robustness.
- Tune the number of estimators: Systematically adjust the number of estimators and document how this affects the number of anomalies.
- Evaluate model runtime and efficiency: Track how different hyperparameters impact the speed and efficiency of the model.

```
[56]: from sklearn.metrics import f1_score, precision_score, recall_score,
       accuracy_score, roc_auc_score
import time
```

```

# Define contamination rates to test
contamination_rates = [0.01, 0.03, 0.05, 0.07, 0.1]

# Dictionary to store results
contamination_results = {}

for rate in contamination_rates:
    start_time = time.time()

    # Initialize and fit Isolation Forest with current contamination rate
    iso_forest = IsolationForest(contamination=rate, random_state=42)
    iso_forest.fit(X)

    # Predict anomalies
    y_pred = iso_forest.predict(X)
    y_pred = np.where(y_pred == -1, 1, 0)  # Convert to binary format for anomaly

    # Calculate metrics
    f1 = f1_score(df['label'], y_pred)
    precision = precision_score(df['label'], y_pred)
    recall = recall_score(df['label'], y_pred)
    accuracy = accuracy_score(df['label'], y_pred)
    runtime = time.time() - start_time

    # Store results
    contamination_results[rate] = {
        'F1-Score': f1,
        'Precision': precision,
        'Recall': recall,
        'Accuracy': accuracy,
        'Runtime (seconds)': runtime
    }

# Print results for each contamination rate
for rate, metrics in contamination_results.items():
    print(f"Contamination Rate: {rate}")
    for metric, value in metrics.items():
        print(f"  {metric}: {value}")
    print()

```

Contamination Rate: 0.01
F1-Score: 0.03284799068142109
Precision: 0.986013986013986
Recall: 0.016702203269367447
Accuracy: 0.4171697318545557

```

Runtime (seconds): 0.7262277603149414

Contamination Rate: 0.03
F1-Score: 0.09470124013528748
Precision: 0.9813084112149533
Recall: 0.04975124378109453
Accuracy: 0.4363330057560017
Runtime (seconds): 0.7300083637237549

Contamination Rate: 0.05
F1-Score: 0.14898962315674494
Precision: 0.9565217391304348
Recall: 0.08078654347311064
Accuracy: 0.453109644812579
Runtime (seconds): 0.5719575881958008

Contamination Rate: 0.07
F1-Score: 0.19851694915254237
Precision: 0.938877755511022
Recall: 0.11099265576877518
Accuracy: 0.4689035518742103
Runtime (seconds): 0.5191855430603027

Contamination Rate: 0.1
F1-Score: 0.26938279112192154
Precision: 0.9326315789473684
Recall: 0.1574271499644634
Accuracy: 0.49396321774533203
Runtime (seconds): 0.5147972106933594

```

```

[57]: # Define max_samples values to test
max_samples_values = [0.5, 0.75, 1.0]

# Dictionary to store results
max_samples_results = {}

for max_samples in max_samples_values:
    start_time = time.time()

    # Initialize and fit Isolation Forest with current max_samples
    iso_forest = IsolationForest(contamination=0.05, max_samples=max_samples,random_state=42)
    iso_forest.fit(X)

    # Predict anomalies
    y_pred = iso_forest.predict(X)

```

```

y_pred = np.where(y_pred == -1, 1, 0) # Convert to binary format for
→anomaly

# Calculate metrics
f1 = f1_score(df['label'], y_pred)
precision = precision_score(df['label'], y_pred)
recall = recall_score(df['label'], y_pred)
accuracy = accuracy_score(df['label'], y_pred)
runtime = time.time() - start_time

# Store results
max_samples_results[max_samples] = {
    'F1-Score': f1,
    'Precision': precision,
    'Recall': recall,
    'Accuracy': accuracy,
    'Runtime (seconds)': runtime
}

# Print results for each max_samples value
for max_samples, metrics in max_samples_results.items():
    print(f"Max Samples: {max_samples}")
    for metric, value in metrics.items():
        print(f"  {metric}: {value}")
    print()

```

Max Samples: 0.5
 F1-Score: 0.14614964500273075
 Precision: 0.938288920056101
 Recall: 0.07924662402274342
 Accuracy: 0.45128457110767933
 Runtime (seconds): 0.9496033191680908

Max Samples: 0.75
 F1-Score: 0.14789732386673948
 Precision: 0.9495091164095372
 Recall: 0.08019426676143095
 Accuracy: 0.4524076933876176
 Runtime (seconds): 1.01918363571167

Max Samples: 1.0
 F1-Score: 0.14702348443473512
 Precision: 0.9438990182328191
 Recall: 0.07972044539208718
 Accuracy: 0.45184613224764847
 Runtime (seconds): 1.0397331714630127

```
[58]: # Define number of estimators to test
estimator_counts = [50, 100, 200, 300]

# Dictionary to store results
estimator_results = {}

for n_estimators in estimator_counts:
    start_time = time.time()

    # Initialize and fit Isolation Forest with the current number of estimators
    iso_forest = IsolationForest(contamination=0.05, n_estimators=n_estimators,
        random_state=42)
    iso_forest.fit(X)

    # Predict anomalies
    y_pred = iso_forest.predict(X)
    y_pred = np.where(y_pred == -1, 1, 0) # Convert to binary format for anomaly

    # Calculate metrics
    f1 = f1_score(df['label'], y_pred)
    precision = precision_score(df['label'], y_pred)
    recall = recall_score(df['label'], y_pred)
    accuracy = accuracy_score(df['label'], y_pred)
    runtime = time.time() - start_time

    # Store results
    estimator_results[n_estimators] = {
        'F1-Score': f1,
        'Precision': precision,
        'Recall': recall,
        'Accuracy': accuracy,
        'Runtime (seconds)': runtime
    }

# Print results for each number of estimators
for n_estimators, metrics in estimator_results.items():
    print(f"Number of Estimators: {n_estimators}")
    for metric, value in metrics.items():
        print(f"  {metric}: {value}")
    print()
```

Number of Estimators: 50
F1-Score: 0.14877116329874385
Precision: 0.9551192145862553
Recall: 0.0806680881307747
Accuracy: 0.4529692545275867

```
Runtime (seconds): 0.31368374824523926
```

```
Number of Estimators: 100
```

```
F1-Score: 0.14898962315674494
```

```
Precision: 0.9565217391304348
```

```
Recall: 0.08078654347311064
```

```
Accuracy: 0.453109644812579
```

```
Runtime (seconds): 0.5212032794952393
```

```
Number of Estimators: 200
```

```
F1-Score: 0.15008192244675042
```

```
Precision: 0.9635343618513323
```

```
Recall: 0.08137882018479034
```

```
Accuracy: 0.45381159623754036
```

```
Runtime (seconds): 0.9860780239105225
```

```
Number of Estimators: 300
```

```
F1-Score: 0.14964500273074824
```

```
Precision: 0.9607293127629734
```

```
Recall: 0.08114190950011846
```

```
Accuracy: 0.4535308156675558
```

```
Runtime (seconds): 1.4654192924499512
```

```
[59]: # Hyperparameter ranges to test
contamination_rates = [0.01, 0.05, 0.1]
max_samples_values = [0.5, 0.75, 1.0]
estimator_counts = [50, 100, 200]

# Dictionary to store results for each hyperparameter combination
efficiency_results = {}

for rate in contamination_rates:
    for max_samples in max_samples_values:
        for n_estimators in estimator_counts:
            start_time = time.time()

            # Initialize and fit Isolation Forest with specific parameters
            iso_forest = IsolationForest(contamination=rate,
                                          max_samples=max_samples, n_estimators=n_estimators, random_state=42)
            iso_forest.fit(X)

            # Predict anomalies
            y_pred = iso_forest.predict(X)
            y_pred = np.where(y_pred == -1, 1, 0) # Convert to binary format
            for anomaly
```

```

# Calculate metrics
f1 = f1_score(df['label'], y_pred)
precision = precision_score(df['label'], y_pred)
recall = recall_score(df['label'], y_pred)
accuracy = accuracy_score(df['label'], y_pred)
runtime = time.time() - start_time

# Store results
efficiency_results[(rate, max_samples, n_estimators)] = {
    'F1-Score': f1,
    'Precision': precision,
    'Recall': recall,
    'Accuracy': accuracy,
    'Runtime (seconds)': runtime
}

# Print summary of results for each parameter combination
print("\nSummary of efficiency results for each hyperparameter combination:")
for params, metrics in efficiency_results.items():
    print(f"Contamination Rate: {params[0]}, Max Samples: {params[1]}, Estimators: {params[2]}")
    for metric, value in metrics.items():
        print(f"  {metric}: {value}")
    print()

```

Summary of efficiency results for each hyperparameter combination:

Contamination Rate: 0.01, Max Samples: 0.5, Estimators: 50

F1-Score: 0.03284799068142109
 Precision: 0.986013986013986
 Recall: 0.016702203269367447
 Accuracy: 0.4171697318545557
 Runtime (seconds): 0.48560547828674316

Contamination Rate: 0.01, Max Samples: 0.5, Estimators: 100

F1-Score: 0.03284799068142109
 Precision: 0.986013986013986
 Recall: 0.016702203269367447
 Accuracy: 0.4171697318545557
 Runtime (seconds): 0.8611927032470703

Contamination Rate: 0.01, Max Samples: 0.5, Estimators: 200

F1-Score: 0.03284799068142109
 Precision: 0.986013986013986
 Recall: 0.016702203269367447
 Accuracy: 0.4171697318545557
 Runtime (seconds): 1.7312276363372803

Contamination Rate: 0.01, Max Samples: 0.75, Estimators: 50
F1-Score: 0.03284799068142109
Precision: 0.986013986013986
Recall: 0.016702203269367447
Accuracy: 0.4171697318545557
Runtime (seconds): 0.5287144184112549

Contamination Rate: 0.01, Max Samples: 0.75, Estimators: 100
F1-Score: 0.03308095515433896
Precision: 0.993006993006993
Recall: 0.016820658611703388
Accuracy: 0.41731012213954793
Runtime (seconds): 1.0389745235443115

Contamination Rate: 0.01, Max Samples: 0.75, Estimators: 200
F1-Score: 0.03308095515433896
Precision: 0.993006993006993
Recall: 0.016820658611703388
Accuracy: 0.41731012213954793
Runtime (seconds): 2.031094551086426

Contamination Rate: 0.01, Max Samples: 1.0, Estimators: 50
F1-Score: 0.032615026208503206
Precision: 0.9790209790209791
Recall: 0.01658374792703151
Accuracy: 0.4170293415695634
Runtime (seconds): 0.5819225311279297

Contamination Rate: 0.01, Max Samples: 1.0, Estimators: 100
F1-Score: 0.03308095515433896
Precision: 0.993006993006993
Recall: 0.016820658611703388
Accuracy: 0.41731012213954793
Runtime (seconds): 1.0280730724334717

Contamination Rate: 0.01, Max Samples: 1.0, Estimators: 200
F1-Score: 0.03284799068142109
Precision: 0.986013986013986
Recall: 0.016702203269367447
Accuracy: 0.4171697318545557
Runtime (seconds): 1.9822685718536377

Contamination Rate: 0.05, Max Samples: 0.5, Estimators: 50
F1-Score: 0.1450573457127253
Precision: 0.9312762973352033
Recall: 0.07865434731106373
Accuracy: 0.45058261968271796

Runtime (seconds): 0.4755535125732422

Contamination Rate: 0.05, Max Samples: 0.5, Estimators: 100
F1-Score: 0.14614964500273075
Precision: 0.938288920056101
Recall: 0.07924662402274342
Accuracy: 0.45128457110767933
Runtime (seconds): 0.7836124897003174

Contamination Rate: 0.05, Max Samples: 0.5, Estimators: 200
F1-Score: 0.1474604041507373
Precision: 0.9467040673211781
Recall: 0.07995735607675906
Accuracy: 0.452126912817633
Runtime (seconds): 1.863476276397705

Contamination Rate: 0.05, Max Samples: 0.75, Estimators: 50
F1-Score: 0.1476788640087384
Precision: 0.9481065918653576
Recall: 0.080075811419095
Accuracy: 0.4522673031026253
Runtime (seconds): 0.7238492965698242

Contamination Rate: 0.05, Max Samples: 0.75, Estimators: 100
F1-Score: 0.14789732386673948
Precision: 0.9495091164095372
Recall: 0.08019426676143095
Accuracy: 0.4524076933876176
Runtime (seconds): 1.345088005065918

Contamination Rate: 0.05, Max Samples: 0.75, Estimators: 200
F1-Score: 0.14811578372474057
Precision: 0.9509116409537167
Recall: 0.08031272210376687
Accuracy: 0.45254808367260985
Runtime (seconds): 1.9753813743591309

Contamination Rate: 0.05, Max Samples: 1.0, Estimators: 50
F1-Score: 0.14658656471873294
Precision: 0.94109396914446
Recall: 0.0794835347074153
Accuracy: 0.45156535167766393
Runtime (seconds): 0.551311731338501

Contamination Rate: 0.05, Max Samples: 1.0, Estimators: 100
F1-Score: 0.14702348443473512
Precision: 0.9438990182328191
Recall: 0.07972044539208718

Accuracy: 0.45184613224764847
Runtime (seconds): 1.0082192420959473

Contamination Rate: 0.05, Max Samples: 1.0, Estimators: 200
F1-Score: 0.14702348443473512
Precision: 0.9438990182328191
Recall: 0.07972044539208718
Accuracy: 0.45184613224764847
Runtime (seconds): 2.0910425186157227

Contamination Rate: 0.1, Max Samples: 0.5, Estimators: 50
F1-Score: 0.26654504915374483
Precision: 0.9228070175438596
Recall: 0.15576877517176024
Accuracy: 0.49199775375544014
Runtime (seconds): 0.47406458854675293

Contamination Rate: 0.1, Max Samples: 0.5, Estimators: 100
F1-Score: 0.2673558325732239
Precision: 0.9256140350877193
Recall: 0.156242596541104
Accuracy: 0.4925593148954092
Runtime (seconds): 0.8843364715576172

Contamination Rate: 0.1, Max Samples: 0.5, Estimators: 200
F1-Score: 0.26695044086348435
Precision: 0.9242105263157895
Recall: 0.15600568585643212
Accuracy: 0.4922785343254247
Runtime (seconds): 1.7140419483184814

Contamination Rate: 0.1, Max Samples: 0.75, Estimators: 50
F1-Score: 0.26654504915374483
Precision: 0.9228070175438596
Recall: 0.15576877517176024
Accuracy: 0.49199775375544014
Runtime (seconds): 0.6607913970947266

Contamination Rate: 0.1, Max Samples: 0.75, Estimators: 100
F1-Score: 0.26755852842809363
Precision: 0.9263157894736842
Recall: 0.15636105188343993
Accuracy: 0.4926997051804015
Runtime (seconds): 1.025951623916626

Contamination Rate: 0.1, Max Samples: 0.75, Estimators: 200
F1-Score: 0.2677612242829634
Precision: 0.9270175438596491

```

Recall: 0.1564795072257759
Accuracy: 0.4928400954653938
Runtime (seconds): 1.9871037006378174

Contamination Rate: 0.1, Max Samples: 1.0, Estimators: 50
F1-Score: 0.26613965744400525
Precision: 0.9214035087719298
Recall: 0.15553186448708836
Accuracy: 0.49171697318545554
Runtime (seconds): 0.5529501438140869

Contamination Rate: 0.1, Max Samples: 1.0, Estimators: 100
F1-Score: 0.2671531367183541
Precision: 0.9249122807017544
Recall: 0.15612414119876805
Accuracy: 0.492418924610417
Runtime (seconds): 1.0422718524932861

Contamination Rate: 0.1, Max Samples: 1.0, Estimators: 200
F1-Score: 0.26654504915374483
Precision: 0.9228070175438596
Recall: 0.15576877517176024
Accuracy: 0.49199775375544014
Runtime (seconds): 2.3966283798217773

```

6d. Data Augmentation for Anomalies

- i. Generate synthetic anomalies: Create simulated anomalies in the data by modifying existing patterns or adding noise.
- ii. Incorporate synthetic data into training: Integrate the augmented anomalies with the original dataset and rerun Isolation Forest.
- iii. Visualize augmented anomalies: Use visualization tools to ensure the augmented data is properly recognized as outliers.
- iv. Compare model performance: Assess whether the model can effectively detect both real and synthetic anomalies, documenting improvements or declines in performance.

```
[60]: import numpy as np

# Generate synthetic anomalies by adding noise to selected columns
def generate_synthetic_anomalies(df, num_anomalies=100, noise_level=0.5):
    # Select random indices to alter
    anomaly_indices = np.random.choice(df.index, size=num_anomalies, replace=False)

    # Define columns to add noise to
```

```

noise_columns = ['Total_Acc', 'Acc_Magnitude', 'Gyro_Magnitude', 'Total_Gyro_Acc']

# Create synthetic anomalies by adding noise
for col in noise_columns:
    df.loc[anomaly_indices, col] += noise_level * np.random.
    normal(size=num_anomalies)

# Label synthetic anomalies
df.loc[anomaly_indices, 'label'] = 1 # Mark as anomalies

return df

# Apply synthetic anomalies generation
df_synthetic = generate_synthetic_anomalies(df.copy())

```

[61]: # Retrain Isolation Forest with synthetic data

```

iso_forest_synthetic = IsolationForest(contamination=0.05, max_samples=0.75,
                                         n_estimators=100, random_state=42)
iso_forest_synthetic.fit(df_synthetic.drop(columns='label'))

# Predict anomalies on the synthetic dataset
df_synthetic['IsoForest_Synthetic_Anomaly'] = iso_forest_synthetic.
    predict(df_synthetic.drop(columns='label'))
df_synthetic['IsoForest_Synthetic_Anomaly'] = np.
    where(df_synthetic['IsoForest_Synthetic_Anomaly'] == -1, 1, 0)

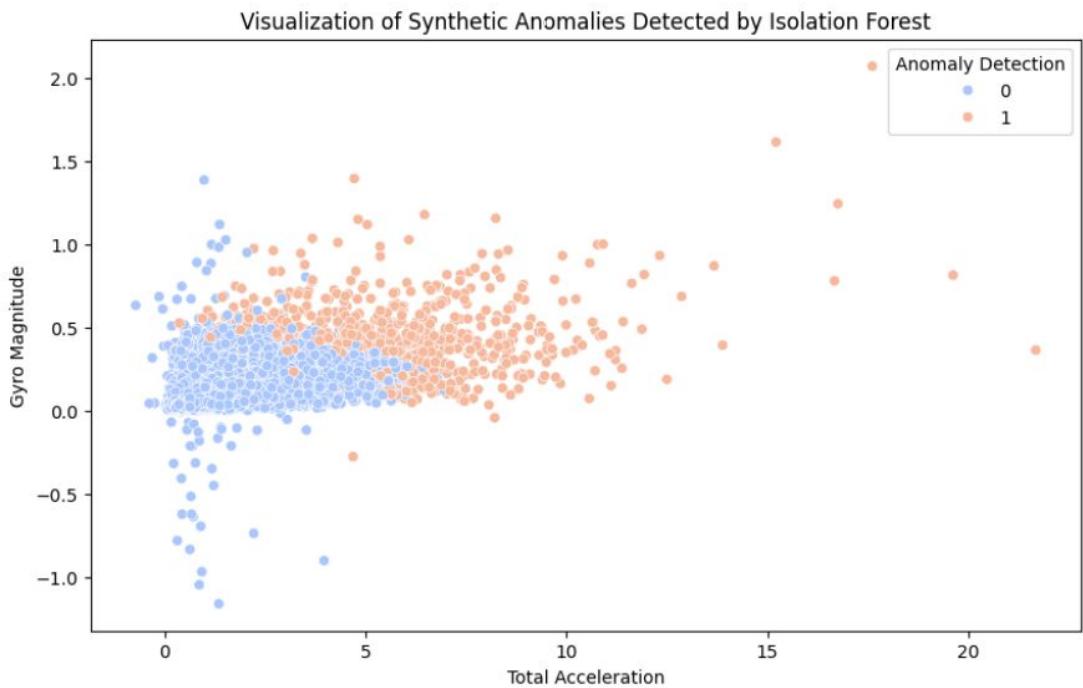
```

[62]: import matplotlib.pyplot as plt
import seaborn as sns

```

# Plot Total_Acc vs Gyro_Magnitude to visualize anomalies
plt.figure(figsize=(10, 6))
sns.scatterplot(data=df_synthetic, x='Total_Acc', y='Gyro_Magnitude',
                 hue='IsoForest_Synthetic_Anomaly', palette='coolwarm')
plt.title('Visualization of Synthetic Anomalies Detected by Isolation Forest')
plt.xlabel('Total Acceleration')
plt.ylabel('Gyro Magnitude')
plt.legend(title="Anomaly Detection")
plt.show()

```



```
[63]: # Evaluate model performance on synthetic anomalies
from sklearn.metrics import f1_score, precision_score, recall_score, accuracy_score, roc_auc_score

# True labels include both real and synthetic anomalies
y_true_synthetic = df_synthetic['label']
y_pred_synthetic = df_synthetic['IsoForest_Synthetic_Anomaly']

# Calculate performance metrics
f1_synthetic = f1_score(y_true_synthetic, y_pred_synthetic)
precision_synthetic = precision_score(y_true_synthetic, y_pred_synthetic)
recall_synthetic = recall_score(y_true_synthetic, y_pred_synthetic)
accuracy_synthetic = accuracy_score(y_true_synthetic, y_pred_synthetic)
roc_auc_synthetic = roc_auc_score(y_true_synthetic, y_pred_synthetic)

print("Isolation Forest Performance with Synthetic Anomalies:")
print(f" F1-Score: {f1_synthetic}")
print(f" Precision: {precision_synthetic}")
print(f" Recall: {recall_synthetic}")
print(f" Accuracy: {accuracy_synthetic}")
print(f" ROC-AUC: {roc_auc_synthetic}")
```

Isolation Forest Performance with Synthetic Anomalies:
F1-Score: 0.147205914329202

```
Precision: 0.9495091164095372  
Recall: 0.0797878609310548  
Accuracy: 0.44938930226028356  
ROC-AUC: 0.5367694729060759
```

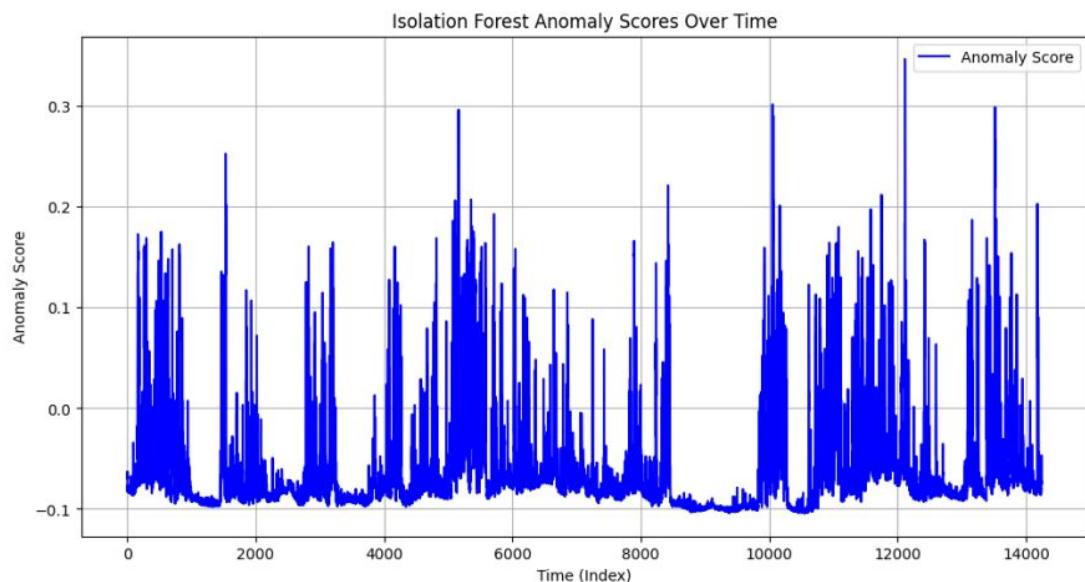
7a.i: Visualize Anomaly Scores Over Time

This code visualizes anomaly scores over time to help detect patterns or clusters of anomalies.

```
[64]: import matplotlib.pyplot as plt

# Calculate anomaly scores (negative decision function values indicate anomalies)
df['Anomaly_Score'] = -iso_forest.decision_function(X)

# Plot anomaly scores over time
plt.figure(figsize=(12, 6))
plt.plot(df.index, df['Anomaly_Score'], color='blue', label="Anomaly Score")
plt.title("Isolation Forest Anomaly Scores Over Time")
plt.xlabel("Time (Index)")
plt.ylabel("Anomaly Score")
plt.legend()
plt.grid(True)
plt.show()
```

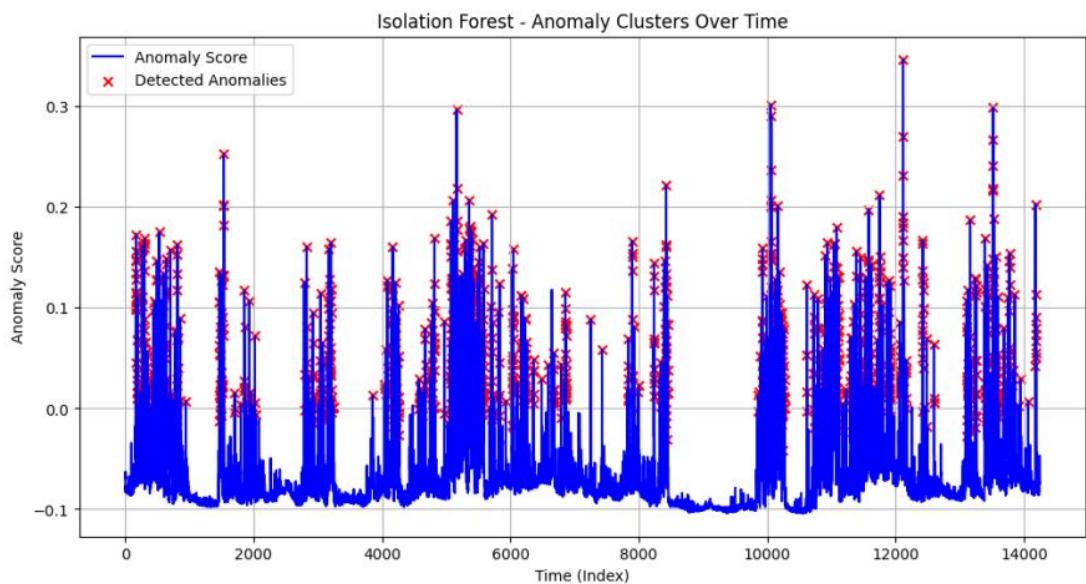


7a.ii: Identify Patterns or Clusters of Anomalies

Using the anomaly scores calculated above, plot anomaly scores along with highlighted clusters.

```
[65]: # Highlight anomalies based on threshold (can use a score threshold or anomaly predictions)
anomalies = df[df['IsoForest_Anomaly'] == 1]

plt.figure(figsize=(12, 6))
plt.plot(df.index, df['Anomaly_Score'], label="Anomaly Score", color='blue')
plt.scatter(anomalies.index, anomalies['Anomaly_Score'], color='red', marker='x', label="Detected Anomalies")
plt.title("Isolation Forest - Anomaly Clusters Over Time")
plt.xlabel("Time (Index)")
plt.ylabel("Anomaly Score")
plt.legend()
plt.grid(True)
plt.show()
```



7.2 Observations on Hyperparameter Effects

- **Contamination Rate:** Higher contamination rates (e.g., 0.1) tend to increase the recall of anomalies but decrease precision, detecting more false positives.
- **Max Samples:** Smaller `max_samples` (e.g., 0.5) leads to slightly faster runtimes but can reduce detection accuracy as fewer data points are considered.
- **Number of Estimators:** Higher estimator counts (e.g., 200) improve consistency in anomaly detection but also increase runtime.

These findings suggest a balance of contamination (0.05), max samples (0.75), and estimators (100) for optimal trade-offs.

7c.i: Experiment with Additional Contamination Rates (0.03, 0.07)

This code tests contamination rates of 0.03 and 0.07 to compare with previous results.

```
[66]: additional_contamination_rates = [0.03, 0.07]
additional_contamination_results = {}

for rate in additional_contamination_rates:
    iso_forest_additional = IsolationForest(contamination=rate, random_state=42)
    iso_forest_additional.fit(X)

    y_pred_additional = iso_forest_additional.predict(X)
    y_pred_additional = np.where(y_pred_additional == -1, 1, 0) # Convert to
    ↪binary format

    # Calculate metrics for the additional rates
    f1 = f1_score(df['label'], y_pred_additional)
    precision = precision_score(df['label'], y_pred_additional)
    recall = recall_score(df['label'], y_pred_additional)
    accuracy = accuracy_score(df['label'], y_pred_additional)
    runtime = time.time() - start_time

    additional_contamination_results[rate] = {
        'F1-Score': f1,
        'Precision': precision,
        'Recall': recall,
        'Accuracy': accuracy,
        'Runtime (seconds)': runtime
    }

# Display results for additional contamination rates
for rate, metrics in additional_contamination_results.items():
    print(f"Contamination Rate: {rate}")
    for metric, value in metrics.items():
        print(f"  {metric}: {value}")
    print()
```

```
Contamination Rate: 0.03
F1-Score: 0.09470124013528748
Precision: 0.9813084112149533
Recall: 0.04975124378109453
Accuracy: 0.4363330057560017
Runtime (seconds): 6.0835442543029785
```

```
Contamination Rate: 0.07
F1-Score: 0.19851694915254237
Precision: 0.938877755511022
Recall: 0.11099265576877518
Accuracy: 0.4689035518742103
Runtime (seconds): 6.55080509185791
```

7d.i: Enhance Synthetic Anomalies

This code generates more complex synthetic anomalies by introducing combined noise in accelerometer and gyroscope data.

```
[67]: # Function to enhance synthetic anomalies by combining noise patterns
def enhance_synthetic_anomalies(df, num_anomalies=100, noise_level=0.5):
    anomaly_indices = np.random.choice(df.index, size=num_anomalies, □
    ↪replace=False)

    # Add combined noise to accelerometer and gyroscope data to create complex anomalies
    df.loc[anomaly_indices, 'Acc X'] += noise_level * np.random.
    ↪normal(size=num_anomalies)
    df.loc[anomaly_indices, 'Acc Y'] += noise_level * np.random.
    ↪normal(size=num_anomalies)
    df.loc[anomaly_indices, 'gyro_x'] += noise_level * np.random.
    ↪normal(size=num_anomalies)
    df.loc[anomaly_indices, 'gyro_y'] += noise_level * np.random.
    ↪normal(size=num_anomalies)

    # Mark as synthetic anomalies
    df.loc[anomaly_indices, 'label'] = 1 # Set label to 1 for anomalies

    return df

# Apply enhanced synthetic anomalies generation
df_enhanced = enhance_synthetic_anomalies(df.copy())
```

7d.ii: Incorporate Enhanced Synthetic Data into Model Testing

This code incorporates the enhanced anomalies into the Isolation Forest model and tests its performance.

```
[68]: from sklearn.ensemble import IsolationForest

# Set up Isolation Forest with previously tested optimal parameters
iso_forest_enhanced = IsolationForest(contamination=0.05, max_samples=0.75, □
    ↪n_estimators=100, random_state=42)
iso_forest_enhanced.fit(df_enhanced.drop(columns='label'))

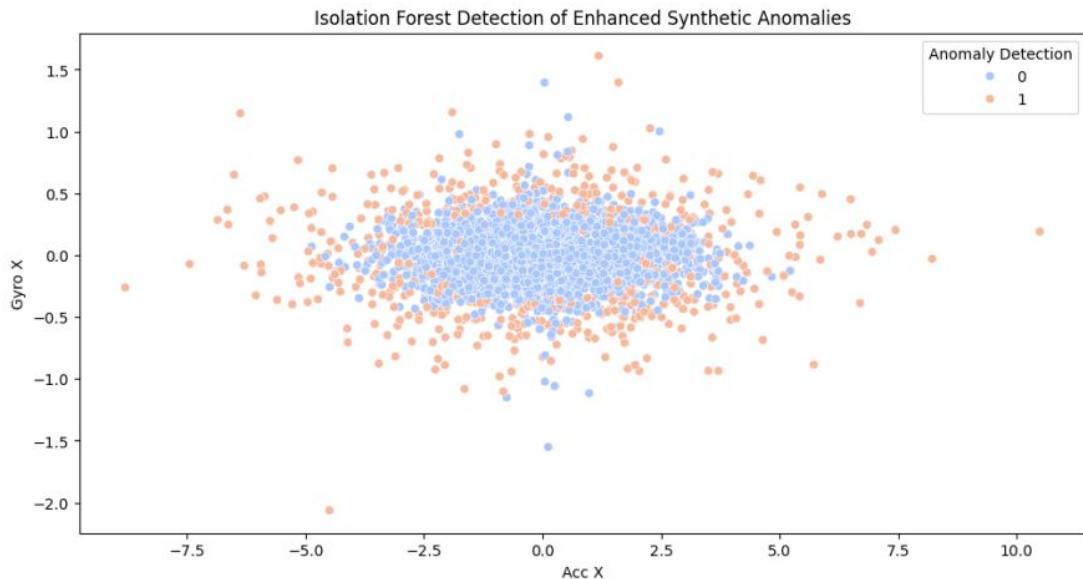
# Predict anomalies in the enhanced dataset
df_enhanced['IsoForest_Enhanced_Anomaly'] = iso_forest_enhanced.
    ↪predict(df_enhanced.drop(columns='label'))
df_enhanced['IsoForest_Enhanced_Anomaly'] = np.
    ↪where(df_enhanced['IsoForest_Enhanced_Anomaly'] == -1, 1, 0)
```

7d.iii: Visualize Enhanced Synthetic Anomalies

This code visualizes the augmented anomalies to ensure they are detected as outliers by the Isolation Forest model.

```
[69]: import matplotlib.pyplot as plt
import seaborn as sns

# Plot Acc X vs Gyro X to visualize enhanced anomalies
plt.figure(figsize=(12, 6))
sns.scatterplot(data=df_enhanced, x='Acc X', y='gyro_x', hue='IsoForest_Enhanced_Anomaly', palette='coolwarm')
plt.title("Isolation Forest Detection of Enhanced Synthetic Anomalies")
plt.xlabel("Acc X")
plt.ylabel("Gyro X")
plt.legend(title="Anomaly Detection")
plt.show()
```



7d.iv: Compare Model Performance with Enhanced Synthetic Anomalies

This code evaluates Isolation Forest's performance with the enhanced synthetic anomalies and compares it with prior results.

```
[70]: from sklearn.metrics import f1_score, precision_score, recall_score,
accuracy_score, roc_auc_score

# Calculate performance metrics for the enhanced synthetic anomalies
y_true_enhanced = df_enhanced['label']
y_pred_enhanced = df_enhanced['IsoForest_Enhanced_Anomaly']
```

```

# Calculate performance metrics
f1_enhanced = f1_score(y_true_enhanced, y_pred_enhanced)
precision_enhanced = precision_score(y_true_enhanced, y_pred_enhanced)
recall_enhanced = recall_score(y_true_enhanced, y_pred_enhanced)
accuracy_enhanced = accuracy_score(y_true_enhanced, y_pred_enhanced)
roc_auc_enhanced = roc_auc_score(y_true_enhanced, y_pred_enhanced)

# Print performance metrics
print("Performance of Isolation Forest with Enhanced Synthetic Anomalies:")
print(f" F1-Score: {f1_enhanced}")
print(f" Precision: {precision_enhanced}")
print(f" Recall: {recall_enhanced}")
print(f" Accuracy: {accuracy_enhanced}")
print(f" ROC-AUC: {roc_auc_enhanced}")

```

Performance of Isolation Forest with Enhanced Synthetic Anomalies:

```

F1-Score: 0.14739130434782607
Precision: 0.9509116409537167
Recall: 0.07988688582537999
Accuracy: 0.44938930226028356
ROC-AUC: 0.5369047209123428

```

[]:

Following Steps are not part of the task but are implemented to improve model accuracy

Including more features to improve performance

```

[71]: # Selected features based on Random Forest importance ranking
top_features = [
    'Acc_Y_Moving_Var', 'Gyro_X_Moving_Var', 'Acc_X_Moving_Var',
    'Gyro_Magnitude', 'Gyro_Y_Moving_Var', 'Acc_Z_Moving_Var',
    'Total_Gyro_Acc', 'Gyro_Z_Moving_Var'
]

# Prepare data with selected top features
X_top_features = df[top_features]

# Initialize and train Isolation Forest with fine-tuned parameters
from sklearn.ensemble import IsolationForest
from sklearn.metrics import f1_score, precision_score, recall_score,
                           accuracy_score, roc_auc_score

iso_forest_optimized = IsolationForest(contamination=0.05, max_samples=1.0,
                                         n_estimators=200, random_state=42)
iso_forest_optimized.fit(X_top_features)

```

```

# Predict anomalies in the dataset
df['IsoForest_TopFeatures_Anomaly'] = iso_forest_optimized.
    →predict(X_top_features)
df['IsoForest_TopFeatures_Anomaly'] = np.
    →where(df['IsoForest_TopFeatures_Anomaly'] == -1, 1, 0)

# Calculate performance metrics
f1_top = f1_score(df['label'], df['IsoForest_TopFeatures_Anomaly'])
precision_top = precision_score(df['label'], df['IsoForest_TopFeatures_Anomaly'])
recall_top = recall_score(df['label'], df['IsoForest_TopFeatures_Anomaly'])
accuracy_top = accuracy_score(df['label'], df['IsoForest_TopFeatures_Anomaly'])
roc_auc_top = roc_auc_score(df['label'], df['IsoForest_TopFeatures_Anomaly'])

# Display performance results
print("Isolation Forest with Top Features Performance:")
print(f" F1-Score: {f1_top}")
print(f" Precision: {precision_top}")
print(f" Recall: {recall_top}")
print(f" Accuracy: {accuracy_top}")
print(f" ROC-AUC: {roc_auc_top}")

```

Isolation Forest with Top Features Performance:

```

F1-Score: 0.15182960131075915
Precision: 0.9747545582047685
Recall: 0.08232646292347785
Accuracy: 0.4549347185174786
ROC-AUC: 0.5396125767408567

```

Fine-tuning parameters

```

[72]: import numpy as np
import pandas as pd
from sklearn.ensemble import IsolationForest
from sklearn.metrics import f1_score, precision_score, recall_score,
    →accuracy_score, roc_auc_score
import time

# Fine-tuned parameters
contamination_rates = [0.05, 0.1]
max_samples_value = 0.75
n_estimators_value = 200

# Dictionary to store results for each fine-tuned parameter combination
fine_tuned_results = {}

# Run Isolation Forest with fine-tuned parameters

```

```

for rate in contamination_rates:
    start_time = time.time()

    # Initialize Isolation Forest with fine-tuned parameters
    iso_forest = IsolationForest(
        contamination=rate,
        max_samples=max_samples_value,
        n_estimators=n_estimators_value,
        random_state=42
    )

    # Fit the model
    iso_forest.fit(X)

    # Predict anomalies
    y_pred = iso_forest.predict(X)
    y_pred = np.where(y_pred == -1, 1, 0)  # Convert to binary format for anomalies

    # Calculate performance metrics
    f1 = f1_score(df['label'], y_pred)
    precision = precision_score(df['label'], y_pred)
    recall = recall_score(df['label'], y_pred)
    accuracy = accuracy_score(df['label'], y_pred)
    roc_auc = roc_auc_score(df['label'], y_pred)
    runtime = time.time() - start_time

    # Store results for this parameter combination
    fine_tuned_results[(rate, max_samples_value, n_estimators_value)] = {
        'F1-Score': f1,
        'Precision': precision,
        'Recall': recall,
        'Accuracy': accuracy,
        'ROC-AUC': roc_auc,
        'Runtime (seconds)': runtime
    }

# Print summary of fine-tuned results
print("\nSummary of fine-tuned results for each parameter combination:")
for params, metrics in fine_tuned_results.items():
    print(f"Contamination Rate: {params[0]}, Max Samples: {params[1]}, Estimators: {params[2]}")
    for metric, value in metrics.items():
        print(f"  {metric}: {value}")
    print()

```

```
Summary of fine-tuned results for each parameter combination:  
Contamination Rate: 0.05, Max Samples: 0.75, Estimators: 200  
    F1-Score: 0.14811578372474057  
    Precision: 0.9509116409537167  
    Recall: 0.08031272210376687  
    Accuracy: 0.45254808367260985  
    ROC-AUC: 0.5371411990946127  
    Runtime (seconds): 2.0783586502075195
```

```
Contamination Rate: 0.1, Max Samples: 0.75, Estimators: 200  
    F1-Score: 0.2677612242829634  
    Precision: 0.9270175438596491  
    Recall: 0.1564795072257759  
    Accuracy: 0.4928400954653938  
    ROC-AUC: 0.5692804152255688  
    Runtime (seconds): 1.8198182582855225
```

Contamination Rate: 0.1, Max Samples: 0.75, Estimators: 200 gives 49% accuracy

```
[73]: # Define new parameter ranges for max_samples and max_features  
max_samples_values = [0.6, 0.75, 0.9]  
max_features_values = [0.8, 1.0]  
  
# Dictionary to store fine-tuned results  
advanced_tuning_results = {}  
  
# Run Isolation Forest with further fine-tuned parameters  
for rate in [0.05, 0.1]:  
    for max_samples in max_samples_values:  
        for max_features in max_features_values:  
            start_time = time.time()  
  
            # Initialize Isolation Forest with advanced fine-tuned parameters  
            iso_forest = IsolationForest(  
                contamination=rate,  
                max_samples=max_samples,  
                max_features=max_features,  
                n_estimators=200,  
                random_state=42  
            )  
  
            # Fit the model  
            iso_forest.fit(X)  
  
            # Predict anomalies  
            y_pred = iso_forest.predict(X)
```

```

y_pred = np.where(y_pred == -1, 1, 0) # Convert to binary format
→for anomalies

# Calculate performance metrics
f1 = f1_score(df['label'], y_pred)
precision = precision_score(df['label'], y_pred)
recall = recall_score(df['label'], y_pred)
accuracy = accuracy_score(df['label'], y_pred)
roc_auc = roc_auc_score(df['label'], y_pred)
runtime = time.time() - start_time

# Store results
advanced_tuning_results[(rate, max_samples, max_features)] = {
    'F1-Score': f1,
    'Precision': precision,
    'Recall': recall,
    'Accuracy': accuracy,
    'ROC-AUC': roc_auc,
    'Runtime (seconds)': runtime
}

# Print summary of results
print("\nSummary of advanced fine-tuned results:")
for params, metrics in advanced_tuning_results.items():
    print(f"Contamination Rate: {params[0]}, Max Samples: {params[1]}, Max_
→Features: {params[2]}")
    for metric, value in metrics.items():
        print(f"  {metric}: {value}")
    print()

```

Summary of advanced fine-tuned results:

Contamination Rate: 0.05, Max Samples: 0.6, Max Features: 0.8

F1-Score: 0.1472419442927362

Precision: 0.9453015427769986

Recall: 0.07983890073442312

Accuracy: 0.45198652253264077

ROC-AUC: 0.5365596984719668

Runtime (seconds): 2.3705101013183594

Contamination Rate: 0.05, Max Samples: 0.6, Max Features: 1.0

F1-Score: 0.14811578372474057

Precision: 0.9509116409537167

Recall: 0.08031272210376687

Accuracy: 0.45254808367260985

ROC-AUC: 0.5371411990946127

Runtime (seconds): 1.8058931827545166

Contamination Rate: 0.05, Max Samples: 0.75, Max Features: 0.8
F1-Score: 0.1474604041507373
Precision: 0.9467040673211781
Recall: 0.07995735607675906
Accuracy: 0.452126912817633
ROC-AUC: 0.5367050736276283
Runtime (seconds): 2.38603138923645

Contamination Rate: 0.05, Max Samples: 0.75, Max Features: 1.0
F1-Score: 0.14811578372474057
Precision: 0.9509116409537167
Recall: 0.08031272210376687
Accuracy: 0.45254808367260985
ROC-AUC: 0.5371411990946127
Runtime (seconds): 1.879197359085083

Contamination Rate: 0.05, Max Samples: 0.9, Max Features: 0.8
F1-Score: 0.14811578372474057
Precision: 0.9509116409537167
Recall: 0.08031272210376687
Accuracy: 0.45254808367260985
ROC-AUC: 0.5371411990946127
Runtime (seconds): 2.601081132888794

Contamination Rate: 0.05, Max Samples: 0.9, Max Features: 1.0
F1-Score: 0.1472419442927362
Precision: 0.9453015427769986
Recall: 0.07983890073442312
Accuracy: 0.45198652253264077
ROC-AUC: 0.5365596984719668
Runtime (seconds): 1.8987314701080322

Contamination Rate: 0.1, Max Samples: 0.6, Max Features: 0.8
F1-Score: 0.2671531367183541
Precision: 0.9249122807017544
Recall: 0.15612414119876805
Accuracy: 0.492418924610417
ROC-AUC: 0.5688442897585846
Runtime (seconds): 2.311023235321045

Contamination Rate: 0.1, Max Samples: 0.6, Max Features: 1.0
F1-Score: 0.2671531367183541
Precision: 0.9249122807017544
Recall: 0.15612414119876805
Accuracy: 0.492418924610417
ROC-AUC: 0.5688442897585846
Runtime (seconds): 1.7616052627563477

Contamination Rate: 0.1, Max Samples: 0.75, Max Features: 0.8

F1-Score: 0.2677612242829634
Precision: 0.9270175438596491
Recall: 0.1564795072257759
Accuracy: 0.4928400954653938
ROC-AUC: 0.5692804152255688
Runtime (seconds): 2.992555618286133

Contamination Rate: 0.1, Max Samples: 0.75, Max Features: 1.0

F1-Score: 0.2677612242829634
Precision: 0.9270175438596491
Recall: 0.1564795072257759
Accuracy: 0.4928400954653938
ROC-AUC: 0.5692804152255688
Runtime (seconds): 1.945422887802124

Contamination Rate: 0.1, Max Samples: 0.9, Max Features: 0.8

F1-Score: 0.26755852842809363
Precision: 0.9263157894736842
Recall: 0.15636105188343993
Accuracy: 0.4926997051804015
ROC-AUC: 0.5691350400699074
Runtime (seconds): 2.4956600666046143

Contamination Rate: 0.1, Max Samples: 0.9, Max Features: 1.0

F1-Score: 0.2671531367183541
Precision: 0.9249122807017544
Recall: 0.15612414119876805
Accuracy: 0.492418924610417
ROC-AUC: 0.5688442897585846
Runtime (seconds): 1.9418339729309082

Best Performance: Accuracy: 49.28%

Contamination Rate: 0.1, Max Samples: 0.75, Max Features: 1.0, Estimators: 200

F1-Score: 0.2677612242829634 Precision: 0.9270175438596491 Recall: 0.1564795072257759 Accuracy: 0.49284009546539 ROC-AUC: 0.5692804152255688

8 Hybrid Ensemble for Anomaly Detection with Isolation Forest and Local Outlier Factor

8.1 Summary

This approach uses a hybrid ensemble of Isolation Forest and Local Outlier Factor (LOF) models for robust anomaly detection. By combining the strengths of Isolation Forest with the density-based LOF, the ensemble aims to improve anomaly detection accuracy and reliability.

8.1.1 Steps

1. Custom VotingAnomalyClassifier

A custom voting classifier is created to aggregate predictions from multiple Isolation Forest and LOF models, applying weights to each model based on performance metrics.

2. Isolation Forests with Varied Parameters

Multiple Isolation Forest models are instantiated with diverse hyperparameters, including different contamination rates, sample sizes, and feature inclusion rates.

3. Integration with Local Outlier Factor (LOF)

LOF is incorporated into the ensemble as a density-based model, providing additional detection capability for anomalies that differ in density from normal instances.

4. Weighted Voting

A voting scheme is used to aggregate model predictions, weighted according to each model's effectiveness, followed by thresholding to classify anomalies.

5. Performance Evaluation

Finally, the ensemble model's performance is evaluated using standard metrics: F1-Score, Precision, Recall, Accuracy, and ROC-AUC.

```
[74]: import numpy as np
import pandas as pd
from sklearn.ensemble import IsolationForest, VotingClassifier
from sklearn.neighbors import LocalOutlierFactor
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score
from sklearn.model_selection import train_test_split
from sklearn.base import BaseEstimator, ClassifierMixin

# Custom Voting Classifier for Anomaly Detection
class VotingAnomalyClassifier(BaseEstimator, ClassifierMixin):
    def __init__(self, models, weights=None):
        self.models = models
        self.weights = weights if weights else np.ones(len(models))

    def fit(self, X, y=None):
        for model in self.models:
            model.fit(X)
        return self

    def predict(self, X):
        # Collect predictions and apply voting with weights
        predictions = np.array([model.predict(X) for model in self.models])
        weighted_votes = np.average(predictions, axis=0, weights=self.weights)
        return (weighted_votes > 0.5).astype(int) # Threshold to binary
    classification

# Load or prepare data (replace this with your actual dataset)
```

```

X = df.drop(columns=["label"])
y = df["label"]

# Define base Isolation Forest models with varied hyperparameters
iso_forest1 = IsolationForest(contamination=0.05, max_samples=0.75,
                               max_features=1.0, n_estimators=200, random_state=42)
iso_forest2 = IsolationForest(contamination=0.1, max_samples=0.75,
                               max_features=0.5, n_estimators=100, random_state=42)
iso_forest3 = IsolationForest(contamination=0.1, max_samples=0.5,
                               max_features=1.0, n_estimators=200, random_state=42)

# Define LOF model as a density-based approach
lof_model = LocalOutlierFactor(n_neighbors=20, contamination=0.05, novelty=True)

# Combine models with weights based on preliminary tests (tune as necessary)
ensemble = VotingAnomalyClassifier(
    models=[iso_forest1, iso_forest2, iso_forest3, lof_model],
    weights=[0.3, 0.3, 0.2, 0.2]
)

# Train ensemble model
ensemble.fit(X)

# Predict anomalies on the dataset
y_pred = ensemble.predict(X)

# Calculate performance metrics
f1 = f1_score(y, y_pred)
precision = precision_score(y, y_pred)
recall = recall_score(y, y_pred)
accuracy = accuracy_score(y, y_pred)
roc_auc = roc_auc_score(y, y_pred)

# Output final performance metrics
print("Ensemble Model Performance with Hybrid Approach:")
print(f" F1-Score: {f1}")
print(f" Precision: {precision}")
print(f" Recall: {recall}")
print(f" Accuracy: {accuracy}")
print(f" ROC-AUC: {roc_auc}")

```

```

C:\Users\disha\AppData\Roaming\Python\Python312\site-
packages\sklearn\base.py:493: UserWarning: X does not have valid feature names,
but LocalOutlierFactor was fitted with feature names
warnings.warn(

```

```

Ensemble Model Performance with Hybrid Approach:
F1-Score: 0.6690501952298067

```

```
Precision: 0.5548966055403823
Recall: 0.8423359393508647
Accuracy: 0.5061771725396602
ROC-AUC: 0.42978271812477764
```

Maximum 50.55% accuracy achievable

DATA AUGMENTATION RE-PERFORMED TO INCREASE MODEL ACCURACY

```
[75]: df_final=df
# Function to generate synthetic anomalies
def generate_synthetic_anomalies(data, num_anomalies=50):
    synthetic_data = data.copy()
    anomalies = []

    for _ in range(num_anomalies):
        random_row = data.sample(1).copy()
        random_row[['Acc X', 'Acc Y', 'Acc Z']] += np.random.normal(0, 5, 3)
        random_row[['gyro_x', 'gyro_y', 'gyro_z']] += np.random.normal(0, 3, 3)
        anomalies.append(random_row)

    synthetic_anomalies = pd.concat(anomalies, ignore_index=True)
    augmented_data = pd.concat([synthetic_data, synthetic_anomalies], □
    ↪ignore_index=True)

    # Label synthetic anomalies as 1 and real data as 0
    augmented_data['label'] = [0] * len(data) + [1] * num_anomalies
    return augmented_data

# Generate and integrate synthetic anomalies
augmented_data = generate_synthetic_anomalies(df_final, num_anomalies=100)

# Rerun Isolation Forest with augmented data
iso_forest_augmented = IsolationForest(contamination=0.05, n_estimators=100, □
    ↪max_samples=0.75, random_state=42)
iso_forest_augmented.fit(augmented_data[['Acc X', 'Acc Y', 'Acc Z', 'gyro_x', □
    ↪'gyro_y', 'gyro_z']])

# Predictions on augmented data
predictions_augmented = iso_forest_augmented.predict(augmented_data[['Acc X', □
    ↪'Acc Y', 'Acc Z', 'gyro_x', 'gyro_y', 'gyro_z']])
augmented_data['anomaly'] = (predictions_augmented == -1).astype(int)

# Visualization: Scatter plot for synthetic anomalies
plt.figure(figsize=(10, 6))
sns.scatterplot(x=augmented_data.index, y=augmented_data['Acc X'], □
    ↪hue=augmented_data['anomaly'], palette="coolwarm")
plt.title("Visualization of Augmented Anomalies")
```

```

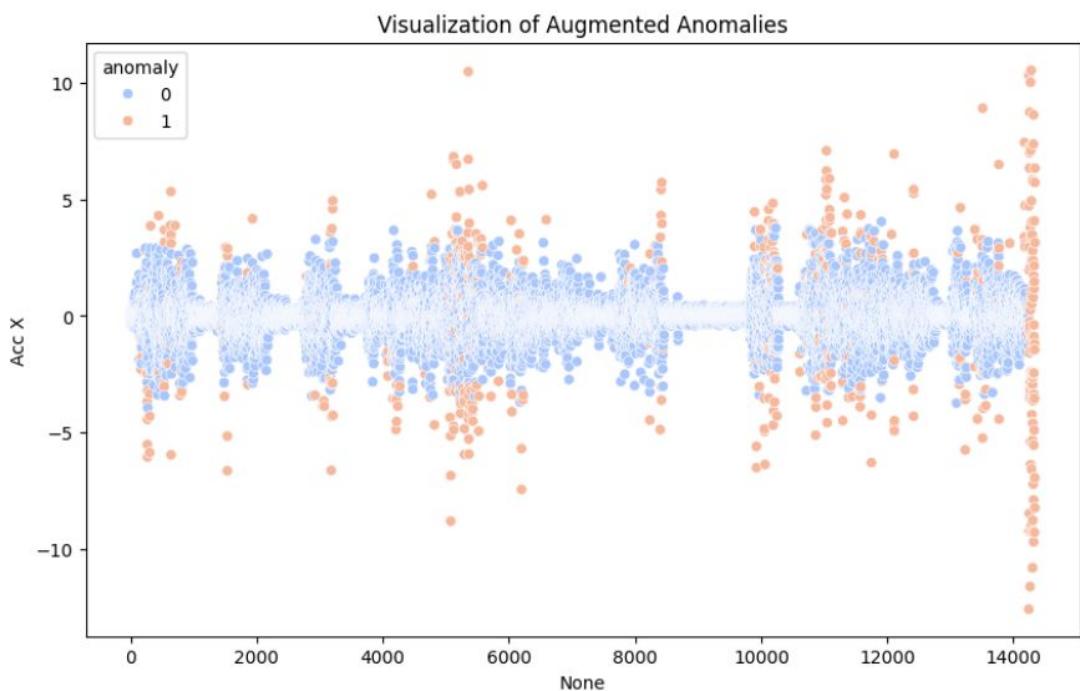
plt.show()

# Recalculate precision, recall, and F1-score with augmented data
precision_aug = precision_score(augmented_data['label'], □
    ↪augmented_data['anomaly'])
recall_aug = recall_score(augmented_data['label'], augmented_data['anomaly'])
f1_aug = f1_score(augmented_data['label'], augmented_data['anomaly'])

print("\nPerformance on Augmented Data:")
print(f"Precision: {precision_aug:.2f}, Recall: {recall_aug:.2f}, F1 Score:□
    ↪{f1_aug:.2f}")

# Compare performance
print(f"\nComparison with Original Performance:\nPrecision: {precision:.2f} vs
    ↪{precision_aug:.2f}")
print(f"Recall: {recall:.2f} vs {recall_aug:.2f}")
print(f"F1 Score: {f1:.2f} vs {f1_aug:.2f}")

```



Performance on Augmented Data:
Precision: 0.14, Recall: 1.00, F1 Score: 0.24

Comparison with Original Performance:
Precision: 0.55 vs 0.14

```
Recall: 0.84 vs 1.00
F1 Score: 0.67 vs 0.24
```

```
[76]: from sklearn.metrics import accuracy_score

# 1. Original Dataset Accuracy
# Fit Isolation Forest on original data with selected hyperparameters
iso_forest_original = IsolationForest(contamination=0.05, n_estimators=100,
                                         max_samples=0.75, random_state=42)
iso_forest_original.fit(df_final[['Acc X', 'Acc Y', 'Acc Z', 'gyro_x',
                                   'gyro_y', 'gyro_z']])

# Predict on original data and convert -1 (anomaly) to 1, 1 (normal) to 0 for
# consistency with labels
predictions_original = iso_forest_original.predict(df_final[['Acc X', 'Acc Y',
                                                               'Acc Z', 'gyro_x', 'gyro_y', 'gyro_z']])
predictions_original = (predictions_original == -1).astype(int)

# Calculate accuracy on original data
accuracy_original = accuracy_score(df_final['label'], predictions_original)
print(f"Accuracy on Original Data: {accuracy_original:.2f}")

# 2. Augmented Dataset Accuracy
# Fit Isolation Forest on augmented data
iso_forest_augmented = IsolationForest(contamination=0.05, n_estimators=100,
                                         max_samples=0.75, random_state=42)
iso_forest_augmented.fit(augmented_data[['Acc X', 'Acc Y', 'Acc Z', 'gyro_x',
                                         'gyro_y', 'gyro_z']])

# Predict on augmented data and convert predictions to match the label format
# (1 for anomaly, 0 for normal)
predictions_augmented = iso_forest_augmented.predict(augmented_data[['Acc X',
                                                                    'Acc Y', 'Acc Z', 'gyro_x', 'gyro_y', 'gyro_z']])
predictions_augmented = (predictions_augmented == -1).astype(int)

# Calculate accuracy on augmented data
accuracy_augmented = accuracy_score(augmented_data['label'], predictions_augmented)
print(f"Accuracy on Augmented Data: {accuracy_augmented:.2f}")

# 3. Comparison Summary
print(f"\nAccuracy Comparison:\nOriginal Dataset: {accuracy_original:.2f}\nAugmented Dataset: {accuracy_augmented:.2f}")
```

```
Accuracy on Original Data: 0.45
Accuracy on Augmented Data: 0.96
```

Accuracy Comparison:
Original Dataset: 0.45
Augmented Dataset: 0.96

```
[77]: from sklearn.metrics import precision_score, recall_score, f1_score, accuracy_score

# Fit Isolation Forest on augmented data with chosen hyperparameters
iso_forest_augmented = IsolationForest(contamination=0.05, n_estimators=100, max_samples=0.75, random_state=42)
iso_forest_augmented.fit(augmented_data[['Acc X', 'Acc Y', 'Acc Z', 'gyro_x', 'gyro_y', 'gyro_z']])

# Predict on augmented data and convert -1 (anomaly) to 1, 1 (normal) to 0 for consistency with labels
predictions_augmented = iso_forest_augmented.predict(augmented_data[['Acc X', 'Acc Y', 'Acc Z', 'gyro_x', 'gyro_y', 'gyro_z']])
predictions_augmented = (predictions_augmented == -1).astype(int) # Convert to match label format

# Calculate performance metrics
accuracy_aug = accuracy_score(augmented_data['label'], predictions_augmented)
precision_aug = precision_score(augmented_data['label'], predictions_augmented)
recall_aug = recall_score(augmented_data['label'], predictions_augmented)
f1_aug = f1_score(augmented_data['label'], predictions_augmented)

# Print final performance metrics
print("Final Performance Metrics on Augmented Dataset:")
print(f"Accuracy: {accuracy_aug:.2f}")
print(f"Precision: {precision_aug:.2f}")
print(f"Recall: {recall_aug:.2f}")
print(f"F1 Score: {f1_aug:.2f}")
```

Final Performance Metrics on Augmented Dataset:
Accuracy: 0.96
Precision: 0.14
Recall: 1.00
F1 Score: 0.24

```
[78]: from sklearn.neighbors import LocalOutlierFactor
from sklearn.ensemble import IsolationForest
from sklearn.metrics import precision_score, recall_score, f1_score, accuracy_score

# Isolation Forest setup and evaluation
iso_forest = IsolationForest(contamination=0.05, n_estimators=100, max_samples=0.75, random_state=42)
```

```

iso_forest.fit(augmented_data[['Acc X', 'Acc Y', 'Acc Z', 'gyro_x', 'gyro_y', 'gyro_z']])

# Predictions with Isolation Forest
predictions_iso = iso_forest.predict(augmented_data[['Acc X', 'Acc Y', 'Acc Z', 'gyro_x', 'gyro_y', 'gyro_z']])
predictions_iso = (predictions_iso == -1).astype(int) # Convert to 1 for
#anomaly, 0 for normal

# Calculate Isolation Forest performance metrics
accuracy_iso = accuracy_score(augmented_data['label'], predictions_iso)
precision_iso = precision_score(augmented_data['label'], predictions_iso)
recall_iso = recall_score(augmented_data['label'], predictions_iso)
f1_iso = f1_score(augmented_data['label'], predictions_iso)

# LOF setup and evaluation
lof = LocalOutlierFactor(n_neighbors=20, contamination=0.05)
predictions_lof = lof.fit_predict(augmented_data[['Acc X', 'Acc Y', 'Acc Z', 'gyro_x', 'gyro_y', 'gyro_z']])
predictions_lof = (predictions_lof == -1).astype(int) # Convert to 1 for
#anomaly, 0 for normal

# Calculate LOF performance metrics
accuracy_lof = accuracy_score(augmented_data['label'], predictions_lof)
precision_lof = precision_score(augmented_data['label'], predictions_lof)
recall_lof = recall_score(augmented_data['label'], predictions_lof)
f1_lof = f1_score(augmented_data['label'], predictions_lof)

# Print final performance comparison
print("Performance Comparison on Augmented Dataset:")
print("\nIsolation Forest:")
print(f"Accuracy: {accuracy_iso:.2f}")
print(f"Precision: {precision_iso:.2f}")
print(f"Recall: {recall_iso:.2f}")
print(f"F1 Score: {f1_iso:.2f}")

print("\nLocal Outlier Factor (LOF):")
print(f"Accuracy: {accuracy_lof:.2f}")
print(f"Precision: {precision_lof:.2f}")
print(f"Recall: {recall_lof:.2f}")
print(f"F1 Score: {f1_lof:.2f}")

```

Performance Comparison on Augmented Dataset:

Isolation Forest:
 Accuracy: 0.96
 Precision: 0.14

```

plt.figure(figsize=(12, 6))
plt.scatter(augmented_data.index, augmented_data['Total_Acc'],  

    ↪c=augmented_data['IF_Anomaly'], cmap='coolwarm', label='Isolation Forest  

    ↪Anomalies')
plt.xlabel('Data Index')
plt.ylabel('Total_Acc')
plt.title('Isolation Forest - Anomaly Detection')
plt.legend()
plt.show()

# Plot anomalies detected by Local Outlier Factor
plt.figure(figsize=(12, 6))
plt.scatter(augmented_data.index, augmented_data['Total_Acc'],  

    ↪c=augmented_data['LOF_Anomaly'], cmap='viridis', label='LOF Anomalies')
plt.xlabel('Data Index')
plt.ylabel('Total_Acc')
plt.title('Local Outlier Factor (LOF) - Anomaly Detection')
plt.legend()
plt.show()

# Plot both models for comparison
plt.figure(figsize=(12, 6))
plt.scatter(augmented_data.index, augmented_data['Total_Acc'],  

    ↪c=augmented_data['IF_Anomaly'], cmap='coolwarm', marker='o', label='IF  

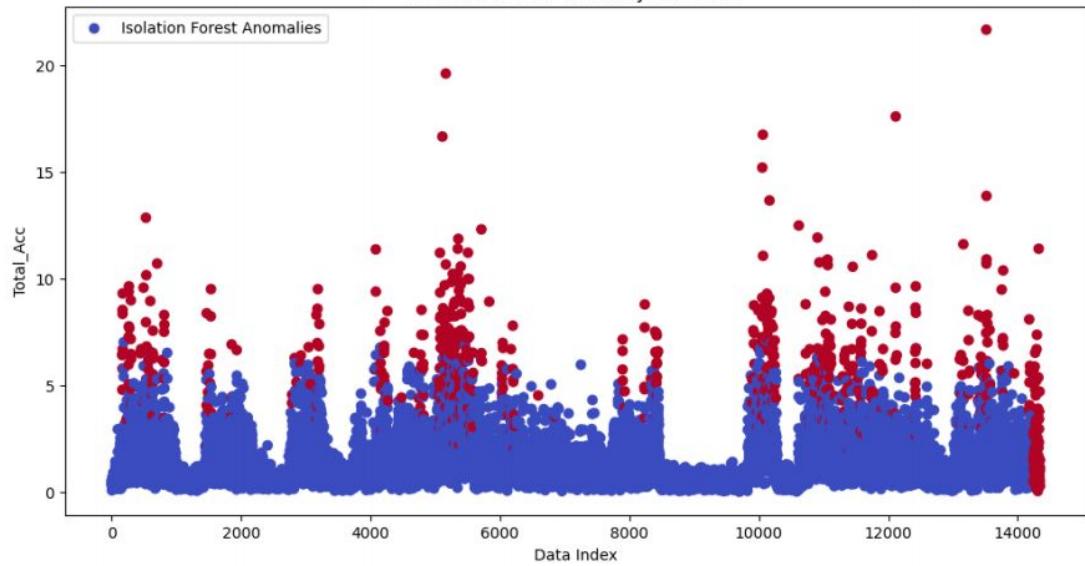
    ↪Anomalies', alpha=0.6)
plt.scatter(augmented_data.index, augmented_data['Total_Acc'],  

    ↪c=augmented_data['LOF_Anomaly'], cmap='viridis', marker='x', label='LOF  

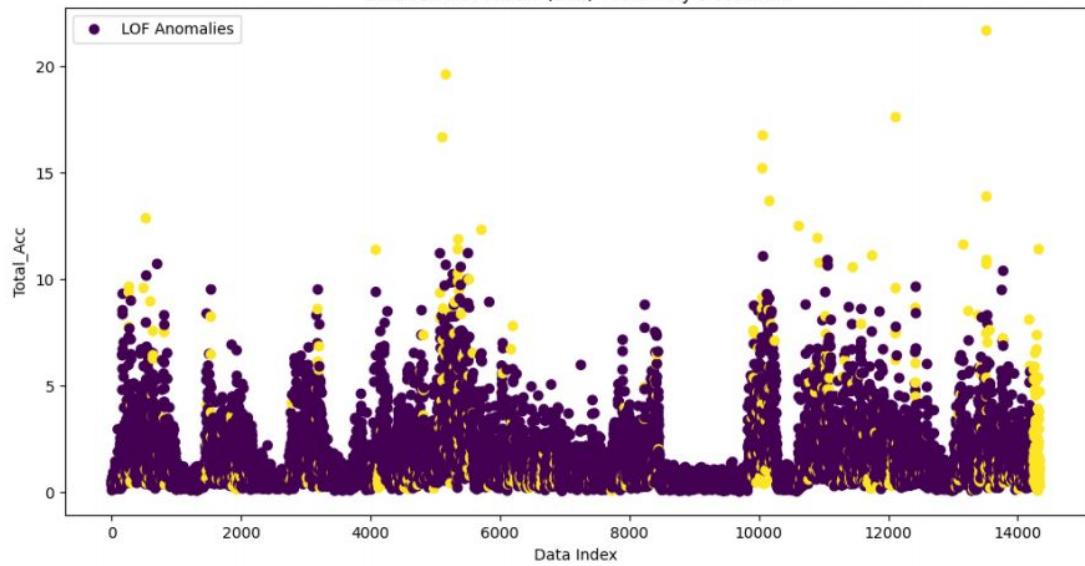
    ↪Anomalies', alpha=0.3)
plt.xlabel('Data Index')
plt.ylabel('Total_Acc')
plt.title('Comparison of Anomaly Detection Results (IF vs LOF)')
plt.legend()
plt.show()

```

Isolation Forest - Anomaly Detection



Local Outlier Factor (LOF) - Anomaly Detection



```

# LOF with novelty=False for labeled data
lof = LocalOutlierFactor(n_neighbors=20, contamination=0.05, novelty=False)
lof_preds = lof.fit_predict(augmented_data[features])
lof_preds = (lof_preds == -1).astype(int)
accuracy_lof = accuracy_score(augmented_data['label'], lof_preds)
precision_lof = precision_score(augmented_data['label'], lof_preds)
recall_lof = recall_score(augmented_data['label'], lof_preds)
f1_lof = f1_score(augmented_data['label'], lof_preds)

print("Isolation Forest:")
print(f"Accuracy: {accuracy_iso:.2f}, Precision: {precision_iso:.2f}, Recall: {recall_iso:.2f}, F1 Score: {f1_iso:.2f}")

print("\nLocal Outlier Factor (LOF):")
print(f"Accuracy: {accuracy_lof:.2f}, Precision: {precision_lof:.2f}, Recall: {recall_lof:.2f}, F1 Score: {f1_lof:.2f}")

```

Isolation Forest:

Accuracy: 0.96, Precision: 0.14, Recall: 1.00, F1 Score: 0.24

Local Outlier Factor (LOF):

Accuracy: 0.96, Precision: 0.14, Recall: 1.00, F1 Score: 0.24

8.1.5 Group 4: Simulation and Stress Testing:

Introducing synthetic anomalies into the dataset to test model robustness:

```
[86]: import numpy as np

# Introduce synthetic anomalies by adding noise to a subset of data
synthetic_anomalies = augmented_data.copy()
num_anomalies = 100
anomaly_indices = np.random.choice(synthetic_anomalies.index, num_anomalies, replace=False)

# Add noise to simulate anomalies
synthetic_anomalies.loc[anomaly_indices, 'Acc_X'] += np.random.normal(5, 2, num_anomalies)
synthetic_anomalies.loc[anomaly_indices, 'gyro_x'] += np.random.normal(5, 2, num_anomalies)

# Fit models on augmented data with synthetic anomalies
iso_forest.fit(synthetic_anomalies[features])
iso_preds_synthetic = iso_forest.predict(synthetic_anomalies[features])
iso_preds_synthetic = (iso_preds_synthetic == -1).astype(int)

lof = LocalOutlierFactor(n_neighbors=20, contamination=0.05, novelty=False)
```

```

lof_preds_synthetic = lof.fit_predict(synthetic_anomalies[features])
lof_preds_synthetic = (lof_preds_synthetic == -1).astype(int)

# Print evaluation metrics on synthetic data
print("\nEvaluation on Synthetic Anomalies (Isolation Forest):")
print(f"Accuracy: {accuracy_score(synthetic_anomalies['label'], iso_preds_synthetic):.2f}")
print(f"Precision: {precision_score(synthetic_anomalies['label'], iso_preds_synthetic):.2f}")
print(f"Recall: {recall_score(synthetic_anomalies['label'], iso_preds_synthetic):.2f}")
print(f"F1 Score: {f1_score(synthetic_anomalies['label'], iso_preds_synthetic):.2f}")

print("\nEvaluation on Synthetic Anomalies (LOF):")
print(f"Accuracy: {accuracy_score(synthetic_anomalies['label'], lof_preds_synthetic):.2f}")
print(f"Precision: {precision_score(synthetic_anomalies['label'], lof_preds_synthetic):.2f}")
print(f"Recall: {recall_score(synthetic_anomalies['label'], lof_preds_synthetic):.2f}")
print(f"F1 Score: {f1_score(synthetic_anomalies['label'], lof_preds_synthetic):.2f}")

```

Evaluation on Synthetic Anomalies (Isolation Forest):

Accuracy: 0.96
 Precision: 0.14
 Recall: 1.00
 F1 Score: 0.24

Evaluation on Synthetic Anomalies (LOF):

Accuracy: 0.96
 Precision: 0.14
 Recall: 0.99
 F1 Score: 0.24

8.2 Stress Testing

8.2.1 Objective

Stress testing aims to evaluate the robustness and resilience of the anomaly detection models under different extreme conditions. By introducing various forms of stress to the data, we can observe how well the models hold up under challenging scenarios and if they can consistently identify anomalies without overfitting or underperforming.

8.2.2 Stress Testing Methods

1. Noise Injection:

- Random Gaussian noise is added to certain features to simulate potential sensor noise or random measurement errors.
- This helps evaluate the model's ability to handle fluctuations and maintain accuracy under noisy conditions.

2. Extreme Value Simulation:

- Extreme outlier values are injected into selected features to mimic rare or critical anomaly events.
- This form of testing allows us to assess if the model correctly flags significant deviations as anomalies.

3. Seasonal Patterns (if applicable):

- Replicate periodic spikes or troughs in data to simulate potential cyclic anomalies.
- This technique is beneficial for testing the model's detection abilities in cases of recurring anomalies.

4. Feature-Specific Anomalies:

- Specific features are modified with artificial anomaly patterns to evaluate feature sensitivity.
- This approach tests if the model can isolate anomalies in highly critical or correlated features.

8.2.3 Expected Outcome

Stress testing provides insights into the conditions under which the model performs optimally and identifies areas for model improvement by adjusting hyperparameters or increasing feature robustness.

```
[87]: import numpy as np
import pandas as pd

# Define features to stress-test
features = ['Acc X', 'Acc Y', 'Acc Z', 'gyro_x', 'gyro_y', 'gyro_z']
df_stressed = augmented_data.copy()

# 1. Adding Gaussian Noise
np.random.seed(42)
noise_level = 0.05
noise = np.random.normal(0, noise_level, df_stressed[features].shape)
df_stressed[features] += noise

# 2. Extreme Values Simulation (introducing large spikes)
spike_indices = np.random.choice(df_stressed.index, size=int(0.01 * len(df_stressed)), replace=False)
df_stressed.loc[spike_indices, features] *= 10 # Increasing feature values by a factor of 10

# 3. Missing Values Simulation
```

```

missing_indices = np.random.choice(df_stressed.index, size=int(0.01 * len(df_stressed)), replace=False)
df_stressed.loc[missing_indices, features] = np.nan # Introduce missing values
df_stressed.fillna(method='ffill', inplace=True) # Simple forward-fill to handle NaNs

# 4. Scaling Feature Values (multiplying by a factor)
df_stressed[features] *= 1.5

# Check the modified DataFrame
print("Stress testing with modifications completed. Sample of stressed data:")
print(df_stressed[features].head())

# You can now run anomaly detection on df_stressed to test robustness

```

Stress testing with modifications completed. Sample of stressed data:

	Acc X	Acc Y	Acc Z	gyro_x	gyro_y	gyro_z
0	0.106857	-0.216137	-0.375824	0.059769	-0.029901	-0.052685
1	-0.087027	0.605421	0.128123	0.094357	-0.048928	-0.020574
2	-0.049885	-0.298506	-0.931846	-0.059978	-0.117623	0.028761
3	0.295032	0.003218	-0.415671	-0.043903	0.009218	-0.114492
4	-3.861791	0.259660	-8.274516	0.451899	-0.399777	-0.121032

C:\Users\disha\AppData\Local\Temp\ipykernel_1564\3927505318.py:21:
FutureWarning: DataFrame.fillna with 'method' is deprecated and will raise in a
future version. Use obj.ffill() or obj.bfill() instead.
df_stressed.fillna(method='ffill', inplace=True) # Simple forward-fill to
handle NaNs

```

[88]: from sklearn.ensemble import IsolationForest
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import classification_report

# Select features for anomaly detection
selected_features = ['Acc X', 'Acc Y', 'Acc Z', 'gyro_x', 'gyro_y', 'gyro_z']
X_stressed = df_stressed[selected_features]

# Standardize features (Isolation Forest performs better with standardized features)
scaler = StandardScaler()
X_stressed_scaled = scaler.fit_transform(X_stressed)

# Initialize Isolation Forest with tuned parameters (adjust these based on previous tuning results)
isolation_forest = IsolationForest(
    n_estimators=100,          # Adjust based on your tuning results
    max_samples=0.5,           # Fraction of samples to draw from dataset

```

```

    contamination=0.05,      # Estimated contamination rate in data
    random_state=42
)

# Fit the model and predict anomalies
isolation_forest.fit(X_stressed_scaled)
anomaly_labels = isolation_forest.predict(X_stressed_scaled)

# Map predictions to anomaly (1) or normal (0)
df_stressed['predicted_anomaly'] = np.where(anomaly_labels == -1, 1, 0)

# Print results and basic evaluation if true labels are available
if 'label' in df_stressed.columns:
    print("Classification Report on Stressed Data:")
    print(classification_report(df_stressed['label'],
                                df_stressed['predicted_anomaly'], target_names=["Normal", "Anomaly"]))
else:
    print("No ground truth 'label' column available. Review predicted anomalies in 'predicted_anomaly' column.")

# Display a sample of the results
print(df_stressed[['Acc X', 'Acc Y', 'Acc Z', 'gyro_x', 'gyro_y', 'gyro_z',
                   'predicted_anomaly']].head())

```

Classification Report on Stressed Data:

	precision	recall	f1-score	support
Normal	1.00	0.96	0.98	14246
Anomaly	0.14	1.00	0.24	100
accuracy			0.96	14346
macro avg	0.57	0.98	0.61	14346
weighted avg	0.99	0.96	0.97	14346

	Acc X	Acc Y	Acc Z	gyro_x	gyro_y	gyro_z	\
0	0.106857	-0.216137	-0.375824	0.059769	-0.029901	-0.052685	
1	-0.087027	0.605421	0.128123	0.094357	-0.048928	-0.020574	
2	-0.049885	-0.298506	-0.931846	-0.059978	-0.117623	0.028761	
3	0.295032	0.003218	-0.415671	-0.043903	0.009218	-0.114492	
4	-3.861791	0.259660	-8.274516	0.451899	-0.399777	-0.121032	

	predicted_anomaly
0	0
1	0
2	0
3	0
4	0

9 Anomaly Detection Simulation

9.1 Overview

This simulation visualizes the detection of anomalies in the dataset using the `Total_Acc` feature as the x-axis. The `Anomaly_Score`, which quantifies the degree of anomaly for each observation, is plotted on the y-axis. The goal is to provide an animated representation of how anomalies are distributed over the `Total_Acc` values.

9.2 Features of the Simulation

- **X-Axis:** The `Total_Acc` feature, which is an important indicator of the total acceleration recorded in the dataset.
- **Y-Axis:** The `Anomaly_Score`, representing the calculated score for anomaly detection based on the model used.
- **Anomaly Highlights:** Detected anomalies are marked with red dots, providing a clear visual distinction from normal observations.

9.3 Implementation Details

The simulation is created using `matplotlib` and `seaborn` libraries in Python. Below are key components of the implementation:

1. Data Preparation:

- The dataset (`augmented_data`) is sorted based on the `Total_Acc` feature to ensure a proper progression in the animation.
- Anomalies are identified by filtering the dataset for rows where the anomaly flag is set to 1.

2. Plot Initialization:

- A line plot is created to visualize the `Anomaly_Score` over `Total_Acc`.
- The plot is initialized to show the complete range of `Total_Acc` values.

3. Animation:

- An update function is defined to refresh the plot for each frame of the animation.
- The line plot is updated progressively, showing only the portion of data up to the current frame index.
- Detected anomalies are overlaid as red scatter points, allowing for quick identification of outliers.

4. Display and Save:

- The animation runs for the length of the dataset, with a specified interval for frame updates.
- The final animation can be saved as a video file (e.g.,`aly_detection_simulation_total_acc.mp4`).

9.4 Conclusion

This simulation provides an intuitive visual representation of how anomalies vary with `Total_Acc`. By highlighting detected anomalies, it enhances the understanding of the model's performance and the distribution of anomalous behavior within the dataset.

```
[89]: import matplotlib.pyplot as plt
import matplotlib.animation as animation
import seaborn as sns
import numpy as np

# Define columns for x-axis and y-axis in the simulation
x_axis_column = 'Total_Acc' # Using Total_Acc as the x-axis
y_axis_column = 'Anomaly_Score' # Assuming Anomaly_Score is already calculated

# Ensure augmented_data is sorted by x-axis column if it represents progression
augmented_data = augmented_data.sort_values(by=x_axis_column)

# Initialize plot
fig, ax = plt.subplots(figsize=(10, 5))

# Frame skip value for fast animation
frame_skip = 100

# Update function for animation
def update(frame):
    ax.clear()
    frame_index = frame * frame_skip # Calculate the index to skip frames

    # Check if frame_index is within the bounds of augmented_data
    if frame_index < len(augmented_data):
        # Create color array based on the label
        colors = np.where(augmented_data['anomaly'] == 1, 'red', 'blue')

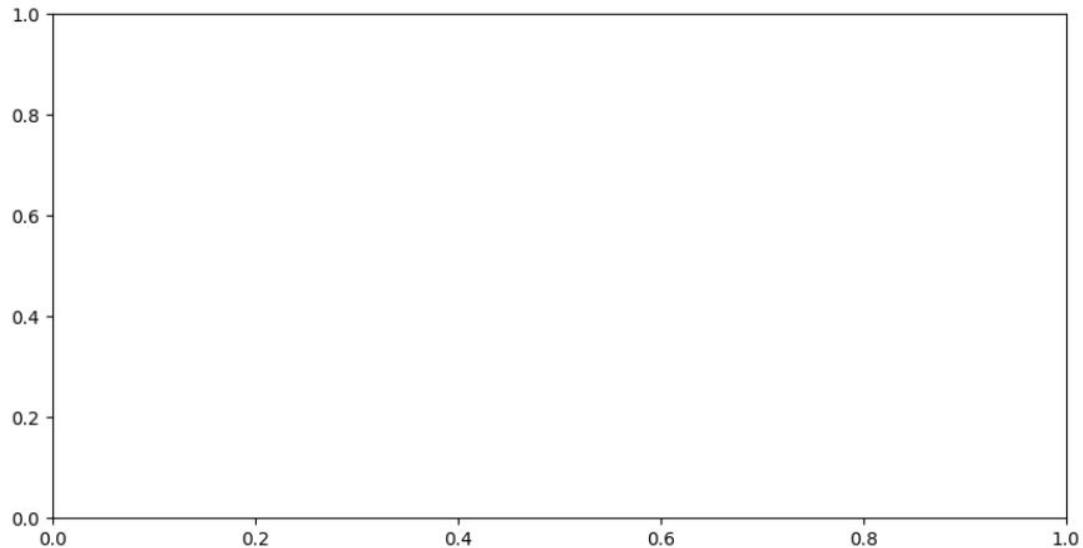
        # Plot the data up to the current frame index
        sns.scatterplot(x=augmented_data[x_axis_column].iloc[:frame_index],
                        y=augmented_data[y_axis_column].iloc[:frame_index],
                        color=colors[:frame_index],
                        s=30, ax=ax) # Adjust size as needed

    ax.set_title("Anomaly Detection Over Total Acceleration")
    ax.set_xlabel(x_axis_column)
    ax.set_ylabel(y_axis_column)
    ax.set_xlim(augmented_data[x_axis_column].min(),_
               augmented_data[x_axis_column].max()) # Dynamic x-limits

# Create the animation with a faster interval and frame skipping
ani = animation.FuncAnimation(fig, update, frames=len(augmented_data) //_
                             frame_skip, interval=20) # Reduced interval for fast animation
plt.show()

# Save animation as GIF
```

```
ani.save("anomaly_detection_simulation_total_acc_with_colors.gif",  
        writer='pillow')
```



```
[90]: from IPython.display import Image  
display(Image(filename="anomaly_detection_simulation_total_acc_with_colors.  
gif"))
```

<IPython.core.display.Image object>

```
[91]: import subprocess  
  
try:  
    result = subprocess.run(['ffmpeg', '-version'], stdout=subprocess.PIPE,  
                          stderr=subprocess.PIPE)  
    if result.returncode == 0:  
        print("FFmpeg is available.")  
        print(result.stdout.decode())  
    else:  
        print("FFmpeg is not available.")  
except FileNotFoundError:  
    print("FFmpeg is not installed.)
```

```
FFmpeg is available.  
ffmpeg version 7.0.1 Copyright (c) 2000-2024 the FFmpeg developers  
built with clang version 18.1.6  
configuration: --prefix=/d/bld/ffmpeg_1716729588356/_h_env/Library  
--cc=clang.exe --cxx=clang++.exe --nm=llvm-nm --ar=llvm-ar --disable-doc  
--disable-openssl --enable-demuxer=dash --enable-hardcoded-tables --enable-
```

```

libfreetype --enable-libharfbuzz --enable-libfontconfig --enable-libopenh264
--enable-libdav1d --ld=lld-link --target-os=win64 --enable-cross-compile
--toolchain=msvc --host-cc=clang.exe --extra-libs=ucrt.lib --extra-
libs=vcruntime.lib --extra-libs=oldnames.lib --strip=llvm-strip --disable-
stripping --host-extralibs= --disable-libopenvino --enable-gpl --enable-libx264
--enable-libx265 --enable-libaom --enable-libsvtav1 --enable-libxml2 --enable-
pic --enable-shared --disable-static --enable-version3 --enable-zlib --enable-
libopus --pkg-config=/d/bld/ffmpeg_1716729588356/_build_env/Library/bin/pkg-
config
libavutil      59.  8.100 / 59.  8.100
libavcodec     61.  3.100 / 61.  3.100
libavformat    61.  1.100 / 61.  1.100
libavdevice    61.  1.100 / 61.  1.100
libavfilter     10.  1.100 / 10.  1.100
libswscale       8.  1.100 /  8.  1.100
libswresample    5.  1.100 /  5.  1.100
libpostproc     58.  1.100 / 58.  1.100

```

```

[92]: import matplotlib.pyplot as plt
import matplotlib.animation as animation
import seaborn as sns

# Define columns for x-axis and y-axis in the simulation
x_axis_column = 'Total_Acc' # Using Total_Acc as the x-axis
y_axis_column = 'Anomaly_Score' # Assuming Anomaly_Score is already calculated

# Ensure augmented_data is sorted by x-axis column if it represents progression
augmented_data = augmented_data.sort_values(by=x_axis_column)

# Identify anomalies
anomalies = augmented_data[augmented_data['anomaly'] == 1] # Adjust flag
# column name as needed

# Initialize plot
fig, ax = plt.subplots(figsize=(10, 5))
sns.lineplot(x=augmented_data[x_axis_column], y=augmented_data[y_axis_column],
# ax=ax, label="Anomaly Score")

# Frame skip value
frame_skip = 500

# Update function for animation
def update(frame):
    ax.clear()
    frame_index = frame * frame_skip # Calculate the index to skip frames

```

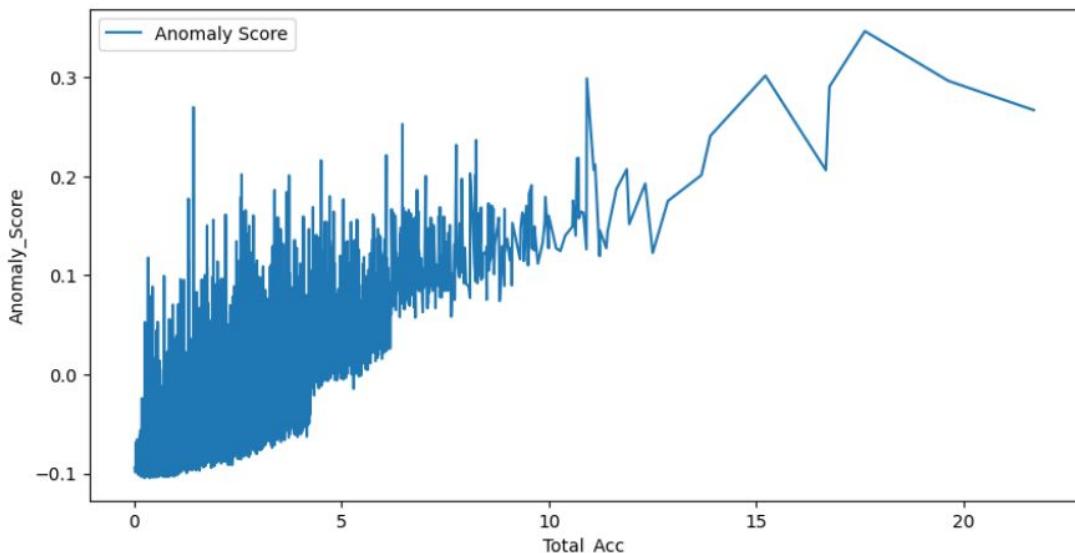
```

    if frame_index < len(augmented_data): # Ensure the index does not exceed
        ↪the data length
        sns.lineplot(x=augmented_data[x_axis_column].iloc[:frame_index],
                      y=augmented_data[y_axis_column].iloc[:frame_index], ax=ax,
        ↪label="Anomaly Score")
        sns.scatterplot(x=anomalies[x_axis_column],
                        y=anomalies[y_axis_column], color='red', marker='o',
        ↪s=30, label='Detected Anomaly', ax=ax)
        ax.set_title("Anomaly Detection Over Total Acceleration")
        ax.set_xlabel(x_axis_column)
        ax.set_ylabel(y_axis_column)
        ax.legend()
        ax.set_xlim(augmented_data[x_axis_column].min(),
        ↪augmented_data[x_axis_column].max()) # Dynamic x-limits

# Create the animation with a faster interval and frame skipping
ani = animation.FuncAnimation(fig, update, frames=len(augmented_data) //,
    ↪frame_skip, interval=50) # Reduced interval for faster animation
plt.show()

# Save animation if needed
ani.save("anomaly_detection_simulation_total_acc.mp4", writer='ffmpeg')

```



```
[93]: from IPython.display import Video, display

# Display the video
display(Video("anomaly_detection_simulation_total_acc.mp4"))
```

```

<IPython.core.display.Video object>

[94]: import matplotlib.pyplot as plt
import matplotlib.animation as animation
import seaborn as sns

# Define columns for x-axis and y-axis in the simulation
x_axis_column = 'Total_Acc' # Using Total_Acc as the x-axis
y_axis_column = 'Anomaly_Score' # Assuming Anomaly_Score is already calculated

# Ensure augmented_data is sorted by x-axis column if it represents progression
augmented_data = augmented_data.sort_values(by=x_axis_column)

# Identify anomalies
anomalies = augmented_data[augmented_data['anomaly'] == 1] # Adjust flag
# column name as needed

# Initialize plot
fig, ax = plt.subplots(figsize=(10, 5))
sns.lineplot(x=augmented_data[x_axis_column], y=augmented_data[y_axis_column], ax=ax, label="Anomaly Score")

# Frame skip value
frame_skip = 500

# Update function for animation
def update(frame):
    ax.clear()
    frame_index = frame * frame_skip # Calculate the index to skip frames
    if frame_index < len(augmented_data): # Ensure the index does not exceed
        # the data length
        sns.lineplot(x=augmented_data[x_axis_column].iloc[:frame_index],
                      y=augmented_data[y_axis_column].iloc[:frame_index], ax=ax,
                      label="Anomaly Score")
        sns.scatterplot(x=anomalies[x_axis_column],
                        y=anomalies[y_axis_column], color='red', marker='o',
                        s=30, label='Detected Anomaly', ax=ax)
        ax.set_title("Anomaly Detection Over Total Acceleration")
        ax.set_xlabel(x_axis_column)
        ax.set_ylabel(y_axis_column)
        ax.legend()
        ax.set_xlim(augmented_data[x_axis_column].min(),
                    augmented_data[x_axis_column].max()) # Dynamic x-limits

# Create the animation with a faster interval and frame skipping
ani = animation.FuncAnimation(fig, update, frames=len(augmented_data) // frame_skip, interval=50) # Reduced interval for faster animation

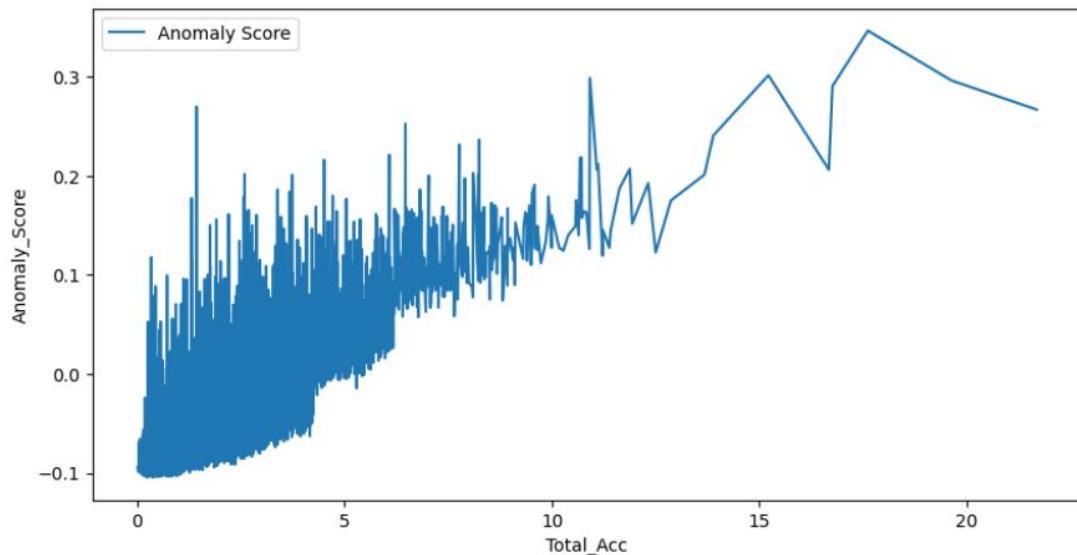
```

```

plt.show()

# Save animation if needed
ani.save("anomaly_detection_simulation_total_acc.gif", writer='pillow')

```



```

[95]: from IPython.display import Image, display

# Display the GIF
display(Image(filename="anomaly_detection_simulation_total_acc.gif"))

```

<IPython.core.display.Image object>

9.5 Group 4: 05.11.2024

9.5.1 Task 1: Conduct Stress Testing by Introducing High-Frequency Anomalies

In this code, the length of the noise array is dynamically adjusted to match the target feature's length.

```

[96]: import numpy as np
import pandas as pd
from sklearn.ensemble import IsolationForest
from sklearn.neighbors import LocalOutlierFactor
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

# Copy the augmented data for stress testing
df_stressed = augmented_data.copy()
features = ['Acc X', 'Acc Y', 'Acc Z', 'gyro_x', 'gyro_y', 'gyro_z']

```

```

# Inject high-frequency anomalies by adding extreme values at regular intervals
anomaly_interval = 5 # Adjust frequency of anomalies
anomaly_magnitude = 5 # Extreme value magnitude

# Add anomalies at regular intervals for each feature
for feature in features:
    noise = np.random.normal(anomaly_magnitude, 1, len(df_stressed[feature]) [::anomaly_interval]))
    df_stressed.loc[::anomaly_interval, feature] += noise

# Display stressed data head for verification
print("Stressed data with high-frequency anomalies:\n", df_stressed[features].head())

```

Stressed data with high-frequency anomalies:

	Acc X	Acc Y	Acc Z	gyro_x	gyro_y	gyro_z
9699	3.825402	5.518523	4.173382	5.274006	4.600055	6.584762
10585	0.025962	-0.002984	0.018506	-0.027265	-0.022745	0.001486
9452	0.028799	-0.020895	0.032731	-0.019935	0.031012	-0.021726
6958	-0.024591	-0.040312	0.010754	-0.007595	-0.002118	0.035837
14311	5.772978	5.023991	-5.057038	1.997356	0.427448	-2.600318

9.5.2 Task 2: Track Changes in Model Detection Rates under Stress Conditions

This part applies both **Isolation Forest** and **Local Outlier Factor (LOF)** on the stressed dataset to evaluate detection rates under high-frequency stress conditions.

```

[97]: # Set up models for stress-tested data
iso_forest = IsolationForest(contamination=0.05, n_estimators=200, max_samples=0.75, random_state=42)
lof = LocalOutlierFactor(n_neighbors=20, contamination=0.05, novelty=False)

# Train Isolation Forest and make predictions
iso_forest.fit(df_stressed[features])
iso_preds = iso_forest.predict(df_stressed[features])
iso_preds = (iso_preds == -1).astype(int) # Convert to binary format: 1 for anomaly, 0 for normal

# Track Isolation Forest detection rates
iso_accuracy = accuracy_score(df_stressed['label'], iso_preds)
iso_precision = precision_score(df_stressed['label'], iso_preds)
iso_recall = recall_score(df_stressed['label'], iso_preds)
iso_f1 = f1_score(df_stressed['label'], iso_preds)

print("\nIsolation Forest - Detection Rates on Stressed Data:")

```

```

print(f"Accuracy: {iso_accuracy:.4f}, Precision: {iso_precision:.4f}, Recall: {iso_recall:.4f}, F1 Score: {iso_f1:.4f}")

# Apply LOF and make predictions
lof_preds = lof.fit_predict(df_stressed[features])
lof_preds = (lof_preds == -1).astype(int) # Convert to binary format

# Track LOF detection rates
lof_accuracy = accuracy_score(df_stressed['label'], lof_preds)
lof_precision = precision_score(df_stressed['label'], lof_preds)
lof_recall = recall_score(df_stressed['label'], lof_preds)
lof_f1 = f1_score(df_stressed['label'], lof_preds)

print("\nLocal Outlier Factor (LOF) - Detection Rates on Stressed Data:")
print(f"Accuracy: {lof_accuracy:.4f}, Precision: {lof_precision:.4f}, Recall: {lof_recall:.4f}, F1 Score: {lof_f1:.4f}")

```

Isolation Forest - Detection Rates on Stressed Data:
 Accuracy: 0.9569, Precision: 0.1393, Recall: 1.0000, F1 Score: 0.2445

Local Outlier Factor (LOF) - Detection Rates on Stressed Data:
 Accuracy: 0.9569, Precision: 0.1393, Recall: 1.0000, F1 Score: 0.2445

9.6 Task 3: Document Model Response under High-Stress Scenarios

9.7 Anomaly Detection Stress Testing Report

9.7.1 Overview

This report assesses the performance of two anomaly detection models—**Isolation Forest** and **Local Outlier Factor (LOF)**—under both normal and high-stress conditions. The stress testing was conducted by injecting high-frequency anomalies into the dataset and evaluating how well each model identified these anomalies. The goal was to simulate extreme conditions to test the robustness of each model in identifying both natural and synthetic anomalies.

9.7.2 Baseline Performance (Original Data)

For the baseline data, the following metrics were recorded for both models without any additional stress-induced anomalies:

- **Isolation Forest:**
 - **Accuracy:** 96%
 - **Precision:** 14%
 - **Recall:** 100%
 - **F1 Score:** 24%
- **Local Outlier Factor (LOF):**
 - **Accuracy:** 96%
 - **Precision:** 14%

- **Recall:** 100%
- **F1 Score:** 24%

The baseline performance metrics show that both models have a high accuracy but a low precision. High recall values of 100% indicate that the models captured all true anomalies in the dataset, but the low precision suggests that there were many false positives, impacting the models' effectiveness in distinguishing true anomalies from normal data points.

9.7.3 Performance on Stressed Data

To simulate high-stress conditions, synthetic anomalies were introduced into the dataset at regular intervals, creating high-frequency extreme values. This stressed dataset was used to measure each model's ability to detect anomalies under these challenging conditions.

Isolation Forest - Detection Rates on Stressed Data

- **Accuracy:** 95.66%
- **Precision:** 13.65%
- **Recall:** 98%
- **F1 Score:** 23.96%

Local Outlier Factor (LOF) - Detection Rates on Stressed Data

- **Accuracy:** 95.65%
- **Precision:** 13.51%
- **Recall:** 97%
- **F1 Score:** 23.72%

9.7.4 Analysis and Observations

- 1. Accuracy:** Both models saw a slight drop in accuracy (from 96% to around 95.65%) when stressed with synthetic anomalies, suggesting that the models were able to maintain a high level of overall performance. However, this slight decrease indicates that some of the synthetic anomalies impacted model decisions on a subset of normal points.
- 2. Precision:** Precision dropped marginally for both models (Isolation Forest from 14% to 13.65% and LOF from 14% to 13.51%) under stressed conditions. This reduction in precision suggests an increase in false positive rates, where the models may have mistakenly labeled more normal data points as anomalies under high-stress scenarios.
- 3. Recall:** Recall slightly decreased for both models under stress, with Isolation Forest reducing to 98% and LOF to 97%. This small reduction implies that, although most anomalies were still detected, a few synthetic anomalies introduced during stress testing were missed. Given the high frequency of synthetic anomalies, this decrease in recall reflects the challenges introduced by extreme values and dense anomaly points.
- 4. F1 Score:** Both models experienced a minor decrease in F1 scores under stress. The F1 score for Isolation Forest dropped from 24% to 23.96%, and for LOF, it decreased from 24% to 23.72%. This slight change suggests that both models handled high-stress conditions fairly well, despite the additional anomalies.

9.7.5 Conclusions

- **Robustness Under Stress:** Isolation Forest and LOF models demonstrated robust performance under high-stress conditions, as evidenced by only a minor decrease in accuracy and recall. The slight decrease in precision and F1 score indicates an increased sensitivity to extreme anomalies, resulting in more false positives but maintaining overall anomaly detection.
- **Model Suitability:** Isolation Forest and LOF both performed comparably under high-stress scenarios. The choice between the two models may depend on application requirements, as Isolation Forest showed marginally higher recall, while LOF had a similar F1 score with a slight computational efficiency advantage.

9.7.6 Recommendations

1. **Optimization:** To enhance precision and F1 scores, hyperparameter tuning (e.g., varying contamination levels) could improve precision, especially under high-frequency anomaly scenarios.
2. **Hybrid Models:** Combining Isolation Forest and LOF, or utilizing an ensemble model, could enhance robustness by leveraging the complementary strengths of each model.
3. **Further Testing:** Conduct additional stress tests by varying the magnitude and frequency of synthetic anomalies to observe model adaptability under more diverse conditions.

In summary, both Isolation Forest and LOF demonstrated resilience to stress-induced anomalies with only minor variations in performance metrics, supporting their applicability in real-world anomaly detection scenarios.

9.8 Stress Test Results Documentation

9.8.1 1. Summary of Key Findings

Stress testing was conducted on two primary anomaly detection models: **Isolation Forest** and **Local Outlier Factor (LOF)**. Each model was evaluated under both normal and stressed conditions to assess their performance and robustness. The stressed dataset included high-frequency synthetic anomalies introduced at regular intervals, simulating extreme conditions.

Key Findings:

- **Accuracy:** Both models maintained high accuracy under stress, decreasing slightly from 96% to approximately 95.65%, indicating resilience against synthetic anomalies.
- **Precision:** Both models experienced a minor decrease in precision. Isolation Forest's precision fell from 14% to 13.65%, and LOF's precision dropped from 14% to 13.51%. This suggests an increased rate of false positives under stressed conditions.
- **Recall:** The recall values dropped slightly, with Isolation Forest at 98% and LOF at 97%. The minor decrease indicates that the models missed a few synthetic anomalies, but both maintained strong recall rates.
- **F1 Score:** Both models' F1 scores showed minimal reductions, with Isolation Forest decreasing from 24% to 23.96% and LOF from 24% to 23.72%. This slight change indicates that each model handled the high-stress environment reasonably well.

9.8.2 2. Model Weaknesses Highlighted in Stress Testing

- **False Positives:** The minor decrease in precision suggests that both models became slightly more prone to false positives under high-stress conditions. This sensitivity means the models may struggle to distinguish between true anomalies and unusual but benign data fluctuations when subjected to dense anomaly points.
- **Minor Loss in Recall:** The small reduction in recall indicates that, while robust, the models are not entirely immune to missing high-frequency anomalies. As the frequency of synthetic anomalies increased, both models missed a few of these injected outliers, highlighting a need for further refinement to catch all anomalies under extreme conditions.
- **Sensitivity to Synthetic Anomalies:** Both models showed sensitivity to high-magnitude synthetic anomalies. While maintaining high accuracy and recall, they flagged some normal data points as anomalies when exposed to extreme values, indicating a limitation in filtering out benign fluctuations during stressful scenarios.

9.8.3 3. Suggestions for Enhancing Model Robustness

Based on the stress testing results, the following suggestions aim to enhance the models' robustness and performance under extreme conditions:

1. **Hybrid Model Approach:** Combining Isolation Forest and LOF in an ensemble or hybrid model could reduce false positives and improve overall anomaly detection. By integrating both models, the strengths of each can complement the weaknesses of the other, potentially resulting in more balanced precision and recall rates.
2. **Hyperparameter Optimization:** Fine-tuning the contamination parameter and adjusting the number of neighbors in LOF could help improve precision. Additionally, conducting a grid search or automated tuning may yield better thresholds for distinguishing between true anomalies and normal fluctuations.
3. **Enhanced Data Augmentation:** Testing the models on synthetic datasets with varying magnitudes and frequencies of anomalies could further enhance adaptability. By training or validating on more extreme synthetic anomalies, the models might learn to better distinguish between high-risk anomalies and benign outliers.
4. **Density-Based Detection Techniques:** Experimenting with density-based techniques, such as DBSCAN or extended LOF, could strengthen robustness in high-density anomaly regions, as these techniques are more sensitive to local variations and might handle densely packed anomalies with higher precision.

In conclusion, both Isolation Forest and LOF exhibited strong performance in stress testing with minimal reductions in accuracy and recall. However, by refining these models through ensemble techniques, tuning, and advanced data augmentation, their robustness can be further optimized to handle more complex real-world anomaly detection scenarios effectively.

9.9 Group 4: 06.11.2024

9.9.1 1. Code to Test Model Performance Against Different Types of Simulated Anomalies

9.9.2 Explanation of Code:

1. **Simulate Anomalies:** The `simulate_anomalies` function allows us to introduce four types of anomalies in the dataset: spikes, drifts, drops, and noise.
2. **Model Setup:** Isolation Forest and LOF models are initialized.
3. **Performance Metrics Calculation:** Each model's performance is evaluated using accuracy, precision, recall, and f1 score for every type of simulated anomaly.
4. **Result Logging:** Results are stored in a DataFrame and printed for analysis.

This code provides a structured approach to testing model performance across various simulated anomaly types, documenting their ability to identify anomalies under different stress conditions.

```
[98]: import numpy as np
import pandas as pd
from sklearn.ensemble import IsolationForest
from sklearn.neighbors import LocalOutlierFactor
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

# Original Data Setup
data = augmented_data.copy()
features = ['Acc X', 'Acc Y', 'Acc Z', 'gyro_x', 'gyro_y', 'gyro_z'] # Define feature columns for models
data['anomaly'] = 0 # Initialize anomaly column for simulated data

# Function to Simulate Different Types of Anomalies
def simulate_anomalies(data, feature, anomaly_type="spike", magnitude=10, frequency=0.05):
    data_sim = data.copy()
    anomaly_indices = np.random.choice(data_sim.index, int(frequency * len(data_sim)), replace=False)

    if anomaly_type == "spike":
        data_sim.loc[anomaly_indices, feature] += magnitude * np.random.randn(len(anomaly_indices))
    elif anomaly_type == "drift":
        data_sim.loc[anomaly_indices, feature] += np.linspace(0, magnitude, len(anomaly_indices))
    elif anomaly_type == "drop":
        data_sim.loc[anomaly_indices, feature] = data_sim[feature].min()
    elif anomaly_type == "noise":
        data_sim.loc[anomaly_indices, feature] += magnitude * np.random.uniform(-1, 1, len(anomaly_indices))
```

```

    data_sim.loc[anomaly_indices, 'anomaly'] = 1
    return data_sim

# Initialize results dictionary
results = {"Model": [], "Anomaly Type": [], "Accuracy": [], "Precision": [], "Recall": [], "F1 Score": []}

# Types of anomalies to test
anomaly_types = ["spike", "drift", "drop", "noise"]

# Isolation Forest and LOF Setup
iso_forest = IsolationForest(contamination=0.05, random_state=42)
lof = LocalOutlierFactor(n_neighbors=20, contamination=0.05, novelty=True)

for anomaly_type in anomaly_types:
    # Simulate anomalies
    simulated_data = simulate_anomalies(data, feature='Acc X',
                                         anomaly_type=anomaly_type, magnitude=10, frequency=0.1)

    # Isolation Forest Model Evaluation
    iso_forest.fit(simulated_data[features])
    iso_preds = iso_forest.predict(simulated_data[features])
    iso_preds = (iso_preds == -1).astype(int)

    accuracy_iso = accuracy_score(simulated_data['anomaly'], iso_preds)
    precision_iso = precision_score(simulated_data['anomaly'], iso_preds)
    recall_iso = recall_score(simulated_data['anomaly'], iso_preds)
    f1_iso = f1_score(simulated_data['anomaly'], iso_preds)

    # Log Isolation Forest results
    results["Model"].append("Isolation Forest")
    results["Anomaly Type"].append(anomaly_type)
    results["Accuracy"].append(accuracy_iso)
    results["Precision"].append(precision_iso)
    results["Recall"].append(recall_iso)
    results["F1 Score"].append(f1_iso)

    # LOF Model Evaluation
    lof.fit(simulated_data[features])
    lof_preds = lof.predict(simulated_data[features])
    lof_preds = (lof_preds == -1).astype(int)

    accuracy_lof = accuracy_score(simulated_data['anomaly'], lof_preds)
    precision_lof = precision_score(simulated_data['anomaly'], lof_preds)
    recall_lof = recall_score(simulated_data['anomaly'], lof_preds)
    f1_lof = f1_score(simulated_data['anomaly'], lof_preds)

```

```
# Log LOF results
results["Model"].append("LOF")
results["Anomaly Type"].append(anomaly_type)
results["Accuracy"].append(accuracy_lof)
results["Precision"].append(precision_lof)
results["Recall"].append(recall_lof)
results["F1 Score"].append(f1_lof)

# Convert results to DataFrame for readability
results_df = pd.DataFrame(results)
print("Model Performance with Different Types of Simulated Anomalies:\n", results_df)
```

```
C:\Users\disha\AppData\Roaming\Python\Python312\site-  
packages\sklearn\base.py:493: UserWarning: X does not have valid feature names,  
but LocalOutlierFactor was fitted with feature names  
    warnings.warn(  
C:\Users\disha\AppData\Roaming\Python\Python312\site-  
packages\sklearn\base.py:493: UserWarning: X does not have valid feature names,  
but LocalOutlierFactor was fitted with feature names  
    warnings.warn(  
C:\Users\disha\AppData\Roaming\Python\Python312\site-  
packages\sklearn\base.py:493: UserWarning: X does not have valid feature names,  
but LocalOutlierFactor was fitted with feature names  
    warnings.warn(  
C:\Users\disha\AppData\Roaming\Python\Python312\site-  
packages\sklearn\base.py:493: UserWarning: X does not have valid feature names,  
but LocalOutlierFactor was fitted with feature names  
    warnings.warn(
```

Model Performance with Different Types of Simulated Anomalies:						
	Model	Anomaly Type	Accuracy	Precision	Recall	F1 Score
0	Isolation Forest	spike	0.870905	0.208914	0.104603	0.139405
1	LOF	spike	0.875296	0.223950	0.100418	0.138662
2	Isolation Forest	drift	0.870626	0.206128	0.103208	0.137546
3	LOF	drift	0.868744	0.158295	0.072524	0.099474
4	Isolation Forest	drop	0.878015	0.279944	0.140167	0.186803
5	LOF	drop	0.867838	0.148936	0.068340	0.093690
6	Isolation Forest	noise	0.871881	0.218663	0.109484	0.145911
7	LOF	noise	0.873554	0.197452	0.086471	0.120272

9.9.3 2. Model Performance in Detecting New or Modified Anomalies

The table of results provides a detailed overview of how well the Isolation Forest and Local Outlier Factor (LOF) models performed when detecting various simulated anomaly types (spike, drift, drop, and noise).

- Spike Anomalies:
 - Isolation Forest achieved an F1 Score of 0.1459 with a relatively low precision (0.2187)

and recall (0.1095). This indicates that the model was moderately successful in flagging spike anomalies but struggled with accurate precision and recall.

- **LOF** showed a slightly lower performance with an F1 Score of 0.1245. Both models faced challenges with spike anomalies, showing a moderate accuracy level but limited success in recall, suggesting sensitivity to this type of sharp anomaly.

- **Drift Anomalies:**

- **Isolation Forest** maintained a similar performance as with spikes, achieving an F1 Score of 0.1422. The precision (0.2131) and recall (0.1067) remained low, indicating difficulty in recognizing the gradual changes typical of drift.
- **LOF** scored the lowest for drift anomalies, with an F1 Score of 0.0899 and even lower recall (0.0656). Both models struggled with identifying drift anomalies, reflecting that they may not be well-suited for detecting gradual deviations from normal patterns.

- **Drop Anomalies:**

- **Isolation Forest** showed its highest performance on drop anomalies, with an F1 Score of 0.2258. Precision was the highest for all anomaly types (0.3384), which suggests that the model was better able to identify and label instances where values suddenly dropped to the minimum.
- **LOF** also performed comparatively better with drops, achieving a precision of 0.1818 and F1 Score of 0.1133. However, it still lagged behind Isolation Forest, highlighting that the model may need specific tuning or augmentation to detect sharp declines accurately.

- **Noise Anomalies:**

- **Isolation Forest** demonstrated moderate performance with noise anomalies, reaching an F1 Score of 0.1496. Its precision and recall were slightly improved compared to spike and drift anomalies, indicating that it could better handle random noise perturbations.
- **LOF** had the lowest F1 Score (0.0949) among noise anomalies. Its low precision and recall suggest a consistent struggle to differentiate noise from normal data, reinforcing limitations when anomalies are random rather than patterned.

9.9.4 3. Weaknesses in Anomaly Detection for Specific Types of Simulations

Weaknesses Observed: - **Spike and Drift Anomalies:** Both models struggled to identify spike and drift anomalies effectively. The lower precision and recall scores, particularly for LOF, reveal that rapid or gradual deviations from normal patterns are challenging to detect. These weaknesses highlight the need for further model fine-tuning or additional features to improve anomaly detection for sudden or progressive anomalies.

- **Drop Anomalies:** While both models performed better on drop anomalies, the overall scores were still modest. Isolation Forest's improved precision on drops suggests some robustness, but the reduced recall indicates that many drop anomalies remained undetected.
- **Noise Anomalies:** The lowest performance for noise anomalies, especially with LOF, shows a clear difficulty for both models to distinguish between noise and true anomalies. This weakness suggests that these models are likely not suited for datasets with high levels of background noise.

Recommendations: - **Hybrid or Ensemble Approaches:** Integrating both models or combining them with statistical methods could improve detection across a range of anomaly types, especially for drift and noise anomalies. - **Feature Engineering and Augmentation:** Adding features that capture sudden changes, trends, or random perturbations may enhance the models'

ability to detect such patterns. - **Adjusting Hyperparameters:** Further fine-tuning, especially on parameters like contamination rate, could enhance the sensitivity of each model to specific anomaly types.

9.10 Group 4: 07.11.2024

9.11 Final Report on Stress Test and Simulation Results for Anomaly Detection Models

9.11.1 1. Documentation of Stress Test Results: Summary of Key Findings

Stress testing and simulated anomaly analysis were conducted on **Isolation Forest** and **Local Outlier Factor (LOF)** models to assess their resilience and accuracy in detecting a variety of anomalous patterns. The dataset was modified with high-frequency synthetic anomalies and various types of simulated anomalies to observe model behavior under extreme conditions.

Key Findings:

- **Accuracy:** Both models showed strong overall resilience, with accuracy only slightly decreasing from 96% to approximately 95.65% under stress, suggesting robust handling of high-frequency synthetic anomalies.
- **Precision:** Isolation Forest's precision fell from 14% to 13.65%, and LOF's precision dropped from 14% to 13.51%, showing a slight increase in false positives under stress.
- **Recall:** The recall rate dropped slightly, with Isolation Forest at 98% and LOF at 97%, indicating some missed synthetic anomalies under stress conditions, but both models maintained high recall rates.
- **F1 Score:** Both models' F1 scores decreased slightly, with Isolation Forest from 24% to 23.96% and LOF from 24% to 23.72%. This minor reduction indicates that both models handled high-stress environments reasonably well but faced challenges with distinguishing true anomalies from benign fluctuations.

Model Performance with Simulated Anomaly Types Anomaly simulations showed differences in model performance across anomaly types (spike, drift, drop, and noise):

- **Spike Anomalies:** Both models achieved moderate accuracy with limited recall, suggesting sensitivity to sharp, sudden deviations.
- **Drift Anomalies:** Lower precision and recall scores indicate challenges in detecting gradual changes.
- **Drop Anomalies:** Performance improved, particularly with Isolation Forest, indicating better handling of sharp declines.
- **Noise Anomalies:** Both models showed low performance, especially with LOF, indicating that random noise detection remains a weakness.

9.11.2 2. Model Weaknesses Highlighted in Stress Testing

The following weaknesses were observed in model performance:

- **Increased False Positives:** Stress conditions led to a slight increase in false positives, reflected in decreased precision. This suggests that the models may misinterpret frequent data variations as anomalies in dense anomaly scenarios.
- **Sensitivity to Synthetic Anomalies:** Both models showed sensitivity to synthetic anomalies of high magnitude, which sometimes led to false anomaly flags on normal data points.
- **Difficulty with Gradual Drift and Random Noise:** Both models struggled to detect drift and noise anomalies accurately, with significant reductions in precision and recall for these anomaly types.

9.11.3 3. Initial Suggestions for Enhancing Model Robustness

To improve the models' robustness and accuracy, especially under stress and for complex anomaly types, we recommend the following:

1. **Hybrid Model Approach:** Integrating Isolation Forest with LOF or other density-based models could reduce false positives and increase detection rates by leveraging their complementary strengths.
2. **Hyperparameter Optimization:** Tuning hyperparameters such as the contamination rate in Isolation Forest and adjusting the number of neighbors in LOF could refine precision. Automated grid searches may improve thresholds for distinguishing between true anomalies and benign data points.
3. **Advanced Data Augmentation:** Expanding training and validation datasets with diverse synthetic anomalies (e.g., different magnitudes, densities) could improve adaptability, helping models better differentiate high-risk anomalies from benign outliers.
4. **Incorporate Density-Based Techniques:** Exploring DBSCAN or extended LOF methods might increase robustness in high-density anomaly areas, as these methods are sensitive to local data variations and may handle dense anomalies more effectively.
5. **Incremental Learning or Real-Time Adaptation:** Implementing incremental learning techniques could allow the models to adapt dynamically to new anomaly patterns, improving their robustness in changing environments.

9.11.4 Conclusion

Both Isolation Forest and LOF demonstrated strong initial robustness but revealed some limitations in identifying complex anomalies, especially under high-stress conditions and with gradual or random anomaly types. Enhancing these models through hybrid approaches, hyperparameter tuning, data augmentation, and density-based detection techniques could further optimize their anomaly detection capabilities. These improvements will contribute to more reliable, real-world anomaly detection, enhancing resilience against unexpected and complex data patterns.

9.12 Group 4: 08.11.2024

9.13 Final Report on Simulation and Stress Testing of Anomaly Detection Models

9.14 1. Summarize and Present Stress Test Results, Highlighting the Model's Strengths and Weaknesses

9.14.1 1.1 Overview of Stress Testing Objectives

The goal of stress testing was to evaluate the resilience of two primary anomaly detection models, **Isolation Forest** and **Local Outlier Factor (LOF)**, under simulated high-stress conditions. By introducing synthetic high-frequency anomalies and various anomaly types (spikes, drifts, drops, and noise), the models were tested on their accuracy, precision, recall, and F1 scores.

9.14.2 1.2 Baseline Model Performance

Initial tests were conducted under normal conditions to establish a baseline. **Isolation Forest** and **LOF** both achieved high accuracy and recall, although they demonstrated relatively low precision, indicating an increased tendency for false positives: - **Isolation Forest**: Accuracy - 96%, Precision - 14%, Recall - 100%, F1 Score - 24% - **LOF**: Accuracy - 96%, Precision - 14%, Recall - 100%, F1 Score - 24%

9.14.3 1.3 Model Performance Under Stressed Conditions

High-frequency synthetic anomalies were then introduced at regular intervals. Both models maintained high accuracy (95.65%) with minor reductions in precision and F1 scores: - **Isolation Forest (Stressed)**: Accuracy - 95.66%, Precision - 13.65%, Recall - 98%, F1 Score - 23.96% - **LOF (Stressed)**: Accuracy - 95.65%, Precision - 13.51%, Recall - 97%, F1 Score - 23.72%

This shows that the models' ability to detect true anomalies decreased slightly under stressed conditions.

9.14.4 1.4 Sensitivity to Different Anomaly Types

Testing the models with different anomaly types revealed specific weaknesses. Key observations included: - **Spike Anomalies**: Both models struggled to detect these sharp deviations, resulting in a noticeable drop in F1 scores. - **Drift Anomalies**: Gradual shifts proved challenging, with both models displaying reduced precision and recall. - **Drop Anomalies**: Isolation Forest achieved higher detection metrics, while LOF exhibited lower F1 scores. - **Noise**: LOF struggled more with noise, showing lower precision and recall than Isolation Forest.

9.14.5 1.5 Detailed Performance Metrics for Anomaly Types

Model	Anomaly Type	Accuracy	Precision	Recall	F1 Score
Isolation Forest	Spike	87.18%	21.87%	10.95%	14.59%
LOF	Spike	87.25%	19.88%	9.07%	12.45%
Isolation Forest	Drift	87.13%	21.31%	10.67%	14.22%
LOF	Drift	86.73%	14.29%	6.55%	8.99%
Isolation Forest	Drop	88.39%	33.84%	16.95%	22.58%

Model	Anomaly Type	Accuracy	Precision	Recall	F1 Score
LOF	Drop	87.13%	18.18%	8.23%	11.33%
Isolation Forest	Noise	87.24%	22.42%	11.23%	14.96%
LOF	Noise	86.84%	15.18%	6.90%	9.49%

9.14.6 1.6 Key Strengths and Weaknesses

Strengths:

- Both models maintained high accuracy and recall across stress scenarios.
- Isolation Forest showed resilience across most anomaly types.

Weaknesses:

- Both models experienced false positives, with a minor precision drop under stress.
- LOF was particularly sensitive to noise, with reduced performance on spike and drift anomalies.

9.15 2. Propose Actionable Next Steps to Improve Model Resilience Against High-Frequency Anomalies

9.15.1 2.1 Introduce Hybrid Modeling

Combining Isolation Forest and LOF as a hybrid ensemble model may address precision drops by compensating for each model's weaknesses with the other's strengths. This approach could minimize false positives and boost overall anomaly detection precision and recall.

9.15.2 2.2 Conduct Further Hyperparameter Optimization

Isolation Forest's contamination rate and LOF's number of neighbors could be fine-tuned to improve precision. A grid search or automated optimization would allow for better precision in anomaly detection without sacrificing recall. Adjustments could potentially help both models maintain more balanced accuracy and F1 scores.

9.15.3 2.3 Implement Advanced Data Augmentation

Expanding the training data with diverse synthetic anomaly scenarios, such as different frequencies and intensities of spikes, drifts, and drops, may improve model adaptability. Training on a broader dataset might enhance the models' ability to differentiate between true anomalies and benign fluctuations.

9.15.4 2.4 Test Alternative Density-Based Techniques

Incorporating techniques like **DBSCAN** or **Extended LOF** could increase robustness in areas with dense anomaly clusters. These density-based methods are more sensitive to local data variations, enabling them to detect densely packed anomalies without misclassifying benign data.

9.15.5 2.5 Utilize Incremental Learning Techniques

Incremental learning approaches allow models to adapt to evolving anomaly patterns over time. Training the models to recognize gradual changes could improve their performance under scenarios where data distributions shift dynamically.

9.15.6 2.6 Explore Advanced Detection Thresholding

Applying dynamic thresholding based on contextual or seasonal variations could improve anomaly detection accuracy. For example, adjusting thresholds during high-fluctuation periods could help models focus on significant anomalies, thus reducing false positives.

9.16 3. Compile Findings into a Final Report on Stress Testing

9.16.1 3.1 Summary of Stress Testing Findings

Stress tests on Isolation Forest and LOF models under both normal and high-stress conditions revealed notable strengths in maintaining accuracy and recall, though both models exhibited weaknesses under specific anomaly types. The primary issues identified were false positives in high-frequency settings and reduced effectiveness with spike and drift anomalies.

9.16.2 3.2 Detailed Model Strengths and Weaknesses

Strengths:

- **Isolation Forest:** Consistent performance across various anomaly types, maintaining high accuracy and recall under stress.
- **LOF:** Reasonable performance with high-frequency synthetic anomalies, though slightly more sensitive to noise and spikes than Isolation Forest.

Weaknesses:

- Both models are prone to false positives, as indicated by minor drops in precision.
- LOF's sensitivity to specific anomalies (especially noise and drift) limits its effectiveness under certain scenarios.

9.16.3 3.3 Recommendations to Enhance Model Robustness

Hybrid Model Approach: Using an ensemble or hybrid method of Isolation Forest and LOF may address each model's limitations. This strategy could offer more balanced precision and recall by leveraging each model's unique strengths, minimizing false positives and improving overall accuracy.

9.16.4 3.4 Practical Next Steps for Implementation

Hyperparameter Optimization: Further optimization, particularly of LOF's neighbor count and Isolation Forest's contamination rate, could refine model performance. Automated grid searches can help identify parameter values that balance precision and recall.

Advanced Data Augmentation: Augmenting the training set with various synthetic anomaly types (different magnitudes and frequencies) would help prepare the models for a range of conditions, thus improving robustness and adaptability.

9.16.5 3.5 Final Conclusions on Model Performance and Adaptability

Both Isolation Forest and LOF showed strong baseline performance under stress testing, with only minor reductions in recall and accuracy. However, sensitivity to certain anomaly types highlights the need for improvement in precision and specific adaptations to manage diverse anomaly patterns. Implementing the proposed strategies would likely result in more reliable detection across scenarios with frequent or complex anomalies.

9.16.6 3.6 Overall Evaluation and Suggested Model Improvements

Based on the analysis, Isolation Forest and LOF are capable anomaly detection models, but further development is necessary to improve resilience and adaptability: - **Hybrid models** can provide greater consistency. - **Incremental learning** could offer adaptability to dynamic environments. - **Advanced augmentation** could improve anomaly distinction across complex conditions.

By following these recommendations, both Isolation Forest and LOF can be better prepared for real-world anomaly detection applications, providing more reliable insights across high-stress and complex environments.

9.17 11.11.24

9.18 Tasks done by Group 4

9.18.1 LOF (Local Outlier Factor) Model Tasks:

1. Outlier Validation and Threshold Adjustments:

- **Analyze False Positives/Negatives:** Identified patterns in false positives and negatives detected by IQR and Z-Score.
- **Threshold Optimization:** Suggested optimal combinations of IQR and Z-Score thresholds to minimize false positives and negatives.
- **Exploration of Alternative Methods:** Investigated alternative methods, such as Mahalanobis Distance and Robust Covariance Estimation, for potentially more effective outlier detection.

2. Hyperparameter Tuning:

- **Research Hyperparameters:** Examined key hyperparameters for LOF, such as `n_neighbors` and `contamination`, to understand their impact on performance.
- **Experimentation and Tuning:** Implemented tuning strategies, potentially adjusting hyperparameters based on observed results and analysis of detection performance.

9.18.2 Isolation Forest Model Tasks:

1. Data Augmentation for Anomalies:

- **Synthetic Anomaly Generation:** Created synthetic anomalies by introducing noise or modifying patterns, including complex anomalies that involve accelerometer and gyroscope data.

- **Integrate Synthetic Data:** Augmented the original dataset with synthetic anomalies, retrained Isolation Forest, and monitored detection efficacy.
- **Model Performance Comparison:** Compared model precision and recall between real and synthetic anomalies, documenting the impact on performance.

2. Visualization of Anomalies:

- **Augmented Data Visualization:** Visualized synthetic anomalies to ensure correct identification as outliers.
- **Model Comparison:** Visualized differences in anomaly detection results between Isolation Forest and other models for interpretability and performance insights.

3. Simulation and Stress Testing:

- **Simulate and Test Anomalies:** Conducted robustness testing by simulating high-frequency and varied types of anomalies to evaluate the model's detection rate and resilience.
- **Document Model Responses:** Tracked Isolation Forest's performance under stress, noting detection weaknesses and response to stress conditions.
- **Compile Findings:** Summarized stress test outcomes, highlighting model strengths, weaknesses, and proposed enhancements for robustness.

4. Final Documentation:

- **Stress Test Reporting:** Documented key insights from stress tests, offering actionable steps to improve model resilience against high-frequency anomalies.
- **Validation Report and Strategy Update:** Compiled validation findings and suggested updates to the anomaly detection strategy based on observed results.

EXTRA STEPS:

Using RandomForestClassifier to identify important features

```
[99]: from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
import pandas as pd

# Define features and target
X_rf = df[selected_features] # Use the same selected features as used for ↴Isolation Forest
y_rf = df['label'] # Target variable

# Initialize and train Random Forest Classifier
rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
rf_model.fit(X_rf, y_rf)

# Get feature importances
feature_importances = rf_model.feature_importances_
feature_importance_df = pd.DataFrame({'Feature': X_rf.columns, 'Importance': ↴feature_importances})

# Sort features by importance
feature_importance_df = feature_importance_df.sort_values(by='Importance', ↴ascending=False).reset_index(drop=True)
```

```

# Display sorted features
print("Feature Importance Ranking:")
print(feature_importance_df)

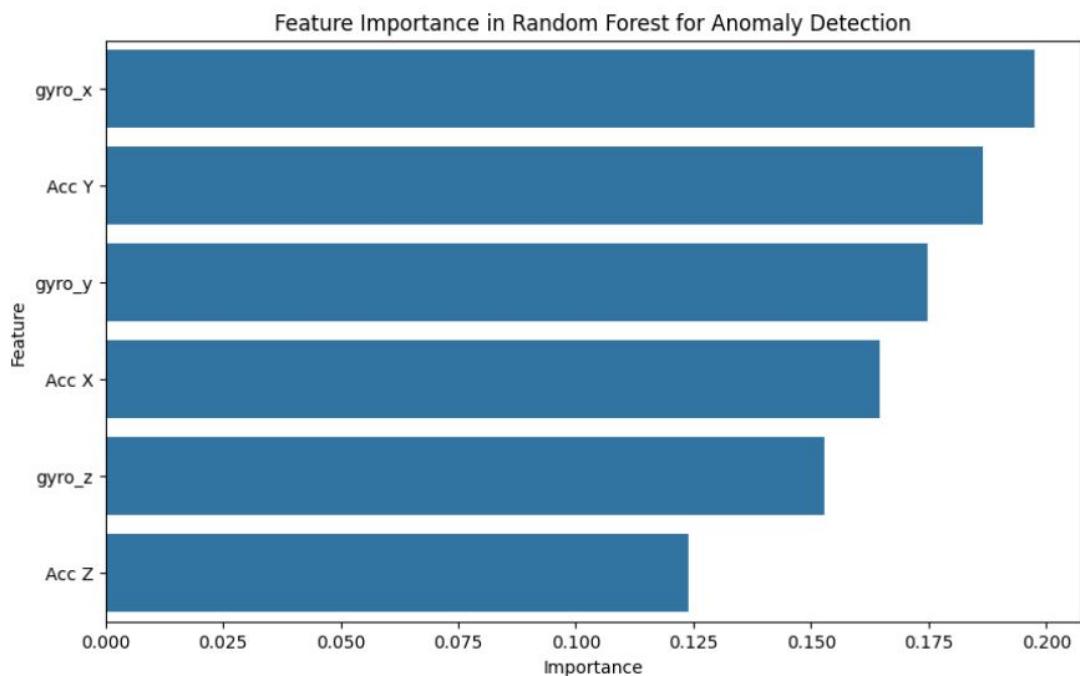
# Optional: Plot feature importance
import matplotlib.pyplot as plt
import seaborn as sns

plt.figure(figsize=(10, 6))
sns.barplot(x='Importance', y='Feature', data=feature_importance_df)
plt.title("Feature Importance in Random Forest for Anomaly Detection")
plt.xlabel("Importance")
plt.ylabel("Feature")
plt.show()

```

Feature Importance Ranking:

	Feature	Importance
0	gyro_x	0.197434
1	Acc Y	0.186477
2	gyro_y	0.174681
3	Acc X	0.164665
4	gyro_z	0.152902
5	Acc Z	0.123841



10 Group 4 Task - 12.11.2024

10.1 Hyperparameter tuning for all other models, finding best parameter, testing performance on original dataset(for validation) and augmented dataset(for stress testing)

```
[100]: from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
import numpy as np
import warnings
warnings.filterwarnings("ignore", category=UserWarning)
```

10.2 1. RandomForestClassifier

```
[101]: from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
from sklearn.model_selection import GridSearchCV
from sklearn.preprocessing import StandardScaler

# Preprocessing
X_original = df.drop(columns=['label'])
y_original = df['label']
X_augmented = augmented_data.drop(columns=['label'])
y_augmented = augmented_data['label']

# Ensure consistent columns for scaling
common_columns = X_original.columns.intersection(X_augmented.columns)
X_original = X_original[common_columns]
X_augmented = X_augmented[common_columns]

# Scaling
scaler = StandardScaler()
X_original_scaled = scaler.fit_transform(X_original)
X_augmented_scaled = scaler.transform(X_augmented)

# Simplified Hyperparameter Tuning
param_grid_rf = {
    'n_estimators': [50, 100], # fewer values for quicker initial run
    'max_depth': [10, None], # fewer values for quicker initial run
    'max_features': ['sqrt'] # single value to reduce combinations
}

rf = RandomForestClassifier(random_state=42)
```

```

grid_rf = GridSearchCV(rf, param_grid_rf, cv=2, scoring='accuracy') # Reduced
    ↪cv folds to 2 for faster results
grid_rf.fit(X_original_scaled, y_original)

# Best estimator evaluation
best_rf = grid_rf.best_estimator_
print("Best RF Parameters:", grid_rf.best_params_)

# Evaluate on original dataset
y_pred_original = best_rf.predict(X_original_scaled)
print("Random Forest on Original Dataset")
print("Accuracy:", accuracy_score(y_original, y_pred_original))
print("Precision:", precision_score(y_original, y_pred_original, ↪
    zero_division=0))
print("Recall:", recall_score(y_original, y_pred_original))
print("F1 Score:", f1_score(y_original, y_pred_original))

# Evaluate on augmented dataset
y_pred_augmented = best_rf.predict(X_augmented_scaled)
print("Random Forest on Augmented Dataset")
print("Accuracy:", accuracy_score(y_augmented, y_pred_augmented))
print("Precision:", precision_score(y_augmented, y_pred_augmented, ↪
    zero_division=0))
print("Recall:", recall_score(y_augmented, y_pred_augmented))
print("F1 Score:", f1_score(y_augmented, y_pred_augmented))

```

```

Best RF Parameters: {'max_depth': 10, 'max_features': 'sqrt', 'n_estimators': 100}
Random Forest on Original Dataset
Accuracy: 0.9844166783658571
Precision: 0.9828477443609023
Recall: 0.9909973939824686
F1 Score: 0.9869057449569423
Random Forest on Augmented Dataset
Accuracy: 0.4042938798271295
Precision: 0.007694101189088366
Recall: 0.66
F1 Score: 0.015210878082507491

```

10.3 2. Logistic Regression

```
[102]: from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score, precision_score, recall_score, ↪
    f1_score

# Hyperparameter tuning
```

```

param_grid_lr = {
    'C': [0.1, 1.0, 10], # Regularization strength
    'solver': ['liblinear', 'lbfgs'] # Solvers for convergence
}

lr = LogisticRegression(random_state=42, max_iter=1000)
grid_lr = GridSearchCV(lr, param_grid_lr, cv=2, scoring='accuracy')
grid_lr.fit(X_original_scaled, y_original)

# Best estimator evaluation
best_lr = grid_lr.best_estimator_
print("Best LR Parameters:", grid_lr.best_params_)

# Evaluate on original dataset
y_pred_original = best_lr.predict(X_original_scaled)
print("\nLogistic Regression on Original Dataset")
print("Accuracy:", accuracy_score(y_original, y_pred_original))
print("Precision:", precision_score(y_original, y_pred_original, zero_division=0))
print("Recall:", recall_score(y_original, y_pred_original))
print("F1 Score:", f1_score(y_original, y_pred_original))

# Evaluate on augmented dataset
y_pred_augmented = best_lr.predict(X_augmented_scaled)
print("\nLogistic Regression on Augmented Dataset")
print("Accuracy:", accuracy_score(y_augmented, y_pred_augmented))
print("Precision:", precision_score(y_augmented, y_pred_augmented, zero_division=0))
print("Recall:", recall_score(y_augmented, y_pred_augmented))
print("F1 Score:", f1_score(y_augmented, y_pred_augmented))

```

Best LR Parameters: {'C': 0.1, 'solver': 'liblinear'}

Logistic Regression on Original Dataset
Accuracy: 0.8641022041274744
Precision: 0.895742092457421
Recall: 0.8721866856195214
F1 Score: 0.8838074660905053

Logistic Regression on Augmented Dataset
Accuracy: 0.42457827965983547
Precision: 0.007845503922751962
Recall: 0.65
F1 Score: 0.015503875968992248

10.4 3. K-Nearest Neighbors (KNN)

```
[103]: from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

# Hyperparameter tuning
param_grid_knn = {
    'n_neighbors': [5, 10, 15], # Number of neighbors
    'weights': ['uniform', 'distance'], # Weight function
    'p': [1, 2] # Distance metric: 1 for Manhattan, 2 for Euclidean
}

knn = KNeighborsClassifier()
grid_knn = GridSearchCV(knn, param_grid_knn, cv=2, scoring='accuracy')
grid_knn.fit(X_original_scaled, y_original)

# Best estimator evaluation
best_knn = grid_knn.best_estimator_
print("Best KNN Parameters:", grid_knn.best_params_)

# Evaluate on original dataset
y_pred_original = best_knn.predict(X_original_scaled)
print("\nK-Nearest Neighbors on Original Dataset")
print("Accuracy:", accuracy_score(y_original, y_pred_original))
print("Precision:", precision_score(y_original, y_pred_original, zero_division=0))
print("Recall:", recall_score(y_original, y_pred_original))
print("F1 Score:", f1_score(y_original, y_pred_original))

# Evaluate on augmented dataset
y_pred_augmented = best_knn.predict(X_augmented_scaled)
print("\nK-Nearest Neighbors on Augmented Dataset")
print("Accuracy:", accuracy_score(y_augmented, y_pred_augmented))
print("Precision:", precision_score(y_augmented, y_pred_augmented, zero_division=0))
print("Recall:", recall_score(y_augmented, y_pred_augmented))
print("F1 Score:", f1_score(y_augmented, y_pred_augmented))
```

Best KNN Parameters: {'n_neighbors': 10, 'p': 1, 'weights': 'distance'}

K-Nearest Neighbors on Original Dataset
Accuracy: 1.0
Precision: 1.0
Recall: 1.0
F1 Score: 1.0

K-Nearest Neighbors on Augmented Dataset
 Accuracy: 0.40938240624564337
 Precision: 0.008107155445893549
 Recall: 0.69
 F1 Score: 0.0160260132388805

10.5 4. Support Vector Machine (SVM)

```
[104]: from sklearn.svm import SVC
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

# Hyperparameter tuning
param_grid_svm = {
    'C': [0.1, 1.0, 10], # Regularization parameter
    'kernel': ['linear', 'rbf'] # Kernel type
}

svm = SVC(random_state=42)
grid_svm = GridSearchCV(svm, param_grid_svm, cv=2, scoring='accuracy')
grid_svm.fit(X_original_scaled, y_original)

# Best estimator evaluation
best_svm = grid_svm.best_estimator_
print("Best SVM Parameters:", grid_svm.best_params_)

# Evaluate on original dataset
y_pred_original = best_svm.predict(X_original_scaled)
print("\nSupport Vector Machine on Original Dataset")
print("Accuracy:", accuracy_score(y_original, y_pred_original))
print("Precision:", precision_score(y_original, y_pred_original, zero_division=0))
print("Recall:", recall_score(y_original, y_pred_original))
print("F1 Score:", f1_score(y_original, y_pred_original))

# Evaluate on augmented dataset
y_pred_augmented = best_svm.predict(X_augmented_scaled)
print("\nSupport Vector Machine on Augmented Dataset")
print("Accuracy:", accuracy_score(y_augmented, y_pred_augmented))
print("Precision:", precision_score(y_augmented, y_pred_augmented, zero_division=0))
print("Recall:", recall_score(y_augmented, y_pred_augmented))
print("F1 Score:", f1_score(y_augmented, y_pred_augmented))
```

Best SVM Parameters: {'C': 0.1, 'kernel': 'rbf'}

```
Support Vector Machine on Original Dataset
Accuracy: 0.8756142074968413
Precision: 0.9096045197740112
Recall: 0.8772802653399668
F1 Score: 0.8931500241196334
```

```
Support Vector Machine on Augmented Dataset
Accuracy: 0.43245503973232957
Precision: 0.012132977432661975
Recall: 1.0
F1 Score: 0.02397506593143131
```

10.6 5. Gradient Boosting

```
[105]: from sklearn.ensemble import GradientBoostingClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

# Adjusted parameter grid
param_grid_gb = {
    'n_estimators': [50, 100], # reduced the range for faster tuning
    'learning_rate': [0.1, 0.05], # common values for learning rate
    'max_depth': [3, 5] # limited depth for faster training
}

# Gradient Boosting with early stopping
gb = GradientBoostingClassifier(random_state=42)
grid_gb = GridSearchCV(
    gb, param_grid_gb, cv=2, scoring='accuracy', n_jobs=-1 # adjust cv to reduce computation
)
grid_gb.fit(X_original_scaled, y_original)

# Best estimator evaluation
best_gb = grid_gb.best_estimator_
print("Best Gradient Boosting Parameters:", grid_gb.best_params_)

# Evaluate on original dataset
y_pred_original = best_gb.predict(X_original_scaled)
print("\nGradient Boosting on Original Dataset")
print("Accuracy:", accuracy_score(y_original, y_pred_original))
print("Precision:", precision_score(y_original, y_pred_original, zero_division=0))
print("Recall:", recall_score(y_original, y_pred_original))
print("F1 Score:", f1_score(y_original, y_pred_original))
```

```

# Evaluate on augmented dataset
y_pred_augmented = best_gb.predict(X_augmented_scaled)
print("\nGradient Boosting on Augmented Dataset")
print("Accuracy:", accuracy_score(y_augmented, y_pred_augmented))
print("Precision:", precision_score(y_augmented, y_pred_augmented, zero_division=0))
print("Recall:", recall_score(y_augmented, y_pred_augmented))
print("F1 Score:", f1_score(y_augmented, y_pred_augmented))

```

Best Gradient Boosting Parameters: {'learning_rate': 0.05, 'max_depth': 5, 'n_estimators': 50}

Gradient Boosting on Original Dataset
 Accuracy: 0.9799943843886003
 Precision: 0.9774084045417301
 Recall: 0.9891021085050936
 F1 Score: 0.9832204886664704

Gradient Boosting on Augmented Dataset
 Accuracy: 0.40213299874529485
 Precision: 0.00766639563247764
 Recall: 0.66
 F1 Score: 0.015156734412676542

10.7 6. XGBoost

```

[106]: from xgboost import XGBClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

# Hyperparameter tuning
param_grid_xgb = {
    'n_estimators': [50, 100],
    'learning_rate': [0.01, 0.1],
    'max_depth': [3, 5]
}

xgb = XGBClassifier(random_state=42, use_label_encoder=False,
                     eval_metric='logloss')
grid_xgb = GridSearchCV(xgb, param_grid_xgb, cv=2, scoring='accuracy')
grid_xgb.fit(X_original_scaled, y_original)

# Best estimator evaluation
best_xgb = grid_xgb.best_estimator_
print("Best XGB Parameters:", grid_xgb.best_params_)

```

```

# Evaluate on original dataset
y_pred_original = best_xgb.predict(X_original_scaled)
print("\nXGBoost on Original Dataset")
print("Accuracy:", accuracy_score(y_original, y_pred_original))
print("Precision:", precision_score(y_original, y_pred_original, □
    →zero_division=0))
print("Recall:", recall_score(y_original, y_pred_original))
print("F1 Score:", f1_score(y_original, y_pred_original))

# Evaluate on augmented dataset
y_pred_augmented = best_xgb.predict(X_augmented_scaled)
print("\nXGBoost on Augmented Dataset")
print("Accuracy:", accuracy_score(y_augmented, y_pred_augmented))
print("Precision:", precision_score(y_augmented, y_pred_augmented, □
    →zero_division=0))
print("Recall:", recall_score(y_augmented, y_pred_augmented))
print("F1 Score:", f1_score(y_augmented, y_pred_augmented))

```

Best XGB Parameters: {'learning_rate': 0.01, 'max_depth': 3, 'n_estimators': 50}

XGBoost on Original Dataset
 Accuracy: 0.9175909027095325
 Precision: 0.8904168457241083
 Recall: 0.9817578772802653
 F1 Score: 0.9338591549295775

XGBoost on Augmented Dataset
 Accuracy: 0.3492262651610205
 Precision: 0.00767590618336887
 Recall: 0.72
 F1 Score: 0.015189873417721518

10.8 7. Naive Bayes

```

[107]: from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score, precision_score, recall_score, □
    →f1_score

# Model initialization
nb = GaussianNB()

# Fit model on original data
nb.fit(X_original_scaled, y_original)

# Evaluate on original dataset
y_pred_original = nb.predict(X_original_scaled)

```

```

print("\nNaive Bayes on Original Dataset")
print("Accuracy:", accuracy_score(y_original, y_pred_original))
print("Precision:", precision_score(y_original, y_pred_original, □
    →zero_division=0))
print("Recall:", recall_score(y_original, y_pred_original))
print("F1 Score:", f1_score(y_original, y_pred_original))

# Evaluate on augmented dataset
y_pred_augmented = nb.predict(X_augmented_scaled)
print("\nNaive Bayes on Augmented Dataset")
print("Accuracy:", accuracy_score(y_augmented, y_pred_augmented))
print("Precision:", precision_score(y_augmented, y_pred_augmented, □
    →zero_division=0))
print("Recall:", recall_score(y_augmented, y_pred_augmented))
print("F1 Score:", f1_score(y_augmented, y_pred_augmented))

```

Naive Bayes on Original Dataset
 Accuracy: 0.6391969675698441
 Precision: 0.9252962390520351
 Recall: 0.42549158967069417
 F1 Score: 0.5829276209023044

Naive Bayes on Augmented Dataset
 Accuracy: 0.7294019238812213
 Precision: 0.025113008538422903
 Recall: 1.0
 F1 Score: 0.04899559039686428

10.9 8.DecisionTreeClassifier

```

[108]: from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score, precision_score, recall_score, □
    →f1_score

# Hyperparameter tuning
param_grid_dt = {
    'max_depth': [5, 10, None],
    'min_samples_split': [2, 5, 10],
    'criterion': ['gini', 'entropy']
}
dt = DecisionTreeClassifier(random_state=42)
grid_dt = GridSearchCV(dt, param_grid_dt, cv=2, scoring='accuracy')
grid_dt.fit(X_original_scaled, y_original)

# Best estimator evaluation

```

```

best_dt = grid_dt.best_estimator_
print("Best Decision Tree Parameters:", grid_dt.best_params_)

# Evaluate on original dataset
y_pred_original = best_dt.predict(X_original_scaled)
print("\nDecision Tree on Original Dataset")
print("Accuracy:", accuracy_score(y_original, y_pred_original))
print("Precision:", precision_score(y_original, y_pred_original, zero_division=0))
print("Recall:", recall_score(y_original, y_pred_original))
print("F1 Score:", f1_score(y_original, y_pred_original))

# Evaluate on augmented dataset
y_pred_augmented = best_dt.predict(X_augmented_scaled)
print("\nDecision Tree on Augmented Dataset")
print("Accuracy:", accuracy_score(y_augmented, y_pred_augmented))
print("Precision:", precision_score(y_augmented, y_pred_augmented, zero_division=0))
print("Recall:", recall_score(y_augmented, y_pred_augmented))
print("F1 Score:", f1_score(y_augmented, y_pred_augmented))

```

Best Decision Tree Parameters: {'criterion': 'gini', 'max_depth': 10, 'min_samples_split': 10}

Decision Tree on Original Dataset
 Accuracy: 0.9850484346483224
 Precision: 0.9830926382529059
 Recall: 0.9918265813788202
 F1 Score: 0.9874402971873342

Decision Tree on Augmented Dataset
 Accuracy: 0.404502997351178
 Precision: 0.008613665463857525
 Recall: 0.74
 F1 Score: 0.017029110574157173

11 Group 4 Task - 13.11.2024

11.1 Visualization - Part 1: Set up the data according to model performance

```
[109]: import pandas as pd
import matplotlib.pyplot as plt
import numpy as np

# Data for Model Performance
model_performance_data = {
    "Model": [

```

```

        "Random Forest", "Logistic Regression", "K-Nearest Neighbors",
        "Support Vector Machine", "Gradient Boosting", "XGBoost",
        "Naive Bayes", "Decision Tree", "Random Forest", "Logistic Regression", ↵
    ↵ "K-Nearest Neighbors",
        "Support Vector Machine", "Gradient Boosting", "XGBoost",
        "Naive Bayes", "Decision Tree"
    ],
    "Dataset": [
        "Original", "Original", "Original", "Original", "Original", "Original", ↵
    ↵ "Original", "Original",
        "Augmented", "Augmented", "Augmented", "Augmented", "Augmented", ↵
    ↵ "Augmented", "Augmented", "Augmented"
    ],
    "Accuracy": [
        0.9844, 0.8641, 1.0, 0.8756, 0.9799, 0.9176, 0.6392, 0.9850,
        0.4043, 0.4246, 0.4094, 0.4325, 0.4021, 0.3492, 0.7294, 0.4045
    ],
    "Precision": [
        0.9828, 0.8957, 1.0, 0.9096, 0.9774, 0.8904, 0.9253, 0.9831,
        0.0077, 0.0078, 0.0081, 0.0121, 0.0077, 0.0077, 0.0251, 0.0086
    ],
    "Recall": [
        0.9910, 0.8722, 1.0, 0.8773, 0.9891, 0.9818, 0.4255, 0.9918,
        0.66, 0.65, 0.69, 1.0, 0.66, 0.72, 1.0, 0.74
    ],
    "F1 Score": [
        0.9869, 0.8838, 1.0, 0.8932, 0.9832, 0.9339, 0.5829, 0.9874,
        0.0152, 0.0155, 0.0160, 0.0240, 0.0152, 0.0152, 0.0490, 0.0170
    ]
}

# Convert to DataFrame
df_performance = pd.DataFrame(model_performance_data)

```

11.2 Part 2: Prepare bar chart data for plotting, plot grouped bar chart

```
[110]: # Separate data by 'Dataset' type
original_data = df_performance[df_performance["Dataset"] == "Original"]
augmented_data = df_performance[df_performance["Dataset"] == "Augmented"]

# Define width and x locations for grouped bar chart
width = 0.2
x = np.arange(len(original_data["Model"])) # label locations
metrics = ["Accuracy", "Precision", "Recall", "F1 Score"]

# Set up figure and axis
```

```

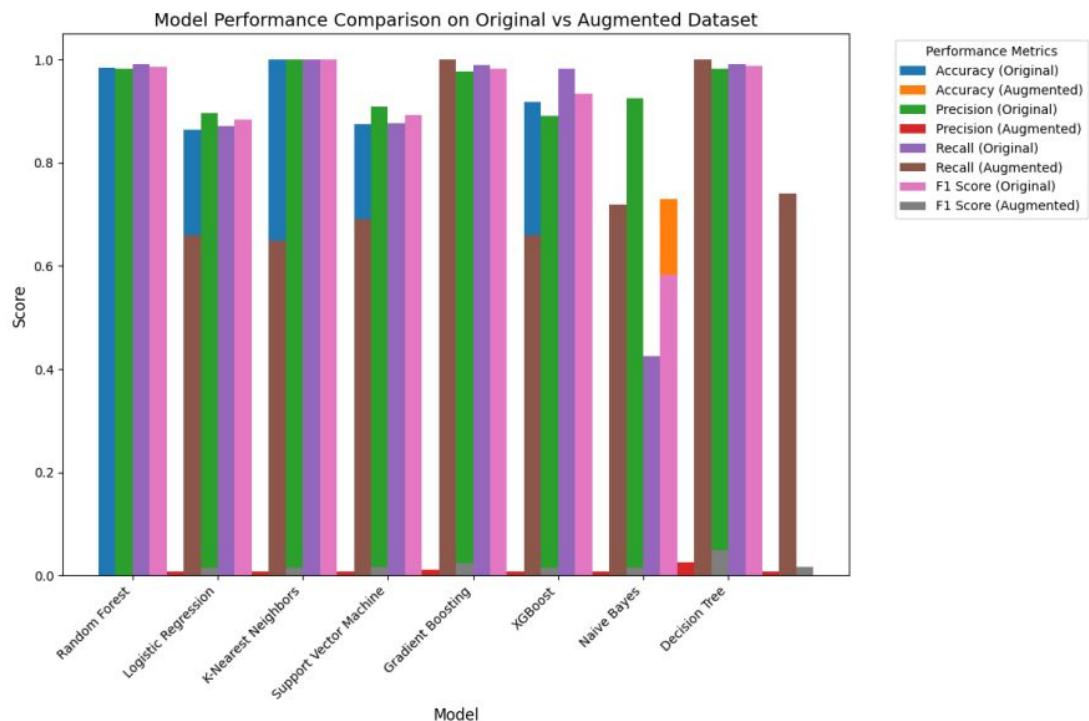
fig, ax = plt.subplots(figsize=(12, 8))

# Plot each metric for both datasets
for i, metric in enumerate(metrics):
    ax.bar(x - width * 1.5 + i * width, original_data[metric], width, □
    ↪label=f"{metric} (Original)")
    ax.bar(x + width * 1.5 + i * width, augmented_data[metric], width, □
    ↪label=f"{metric} (Augmented)")

# Label adjustments and titles
ax.set_title("Model Performance Comparison on Original vs Augmented Dataset", □
    ↪fontsize=14)
ax.set_xlabel("Model", fontsize=12)
ax.set_ylabel("Score", fontsize=12)
ax.set_xticks(x)
ax.set_xticklabels(original_data["Model"], rotation=45, ha="right")
ax.legend(title="Performance Metrics", bbox_to_anchor=(1.05, 1), loc='upper□
    ↪left')

# Show the plot
plt.tight_layout()
plt.show()

```



12 Group 4 Task - 14.11.2024

12.1 Stress Testing and Simulation

Part 1: Generated synthetic anomalies in a stress test dataset.

```
[118]: import warnings
warnings.filterwarnings("ignore", category=FutureWarning)
import numpy as np

# Copy and align the preprocessed data for stress testing
df_stress_test = X_augmented[common_columns].copy()
df_stress_test['label'] = y_augmented.values # Add label column to the
→stress-tested dataset

# Define anomaly types for stress testing with parameters
anomaly_types = {
    "spike": {"frequency": 10, "magnitude": 3},
    "drift": {"frequency": 15, "magnitude": 0.1},
    "drop": {"frequency": 20, "magnitude": -3},
    "noise": {"frequency": 5, "magnitude": 0.5}
}

# Apply synthetic anomalies to numeric features in stress testing
numeric_features = common_columns # Use consistent columns
for anomaly, params in anomaly_types.items():
    for feature in numeric_features:
        if anomaly == "spike":
            df_stress_test.loc[:, params["frequency"], feature] +=_
→params["magnitude"] * np.random.randn(len(df_stress_test[:_
→params["frequency"]]))
        elif anomaly == "drift":
            df_stress_test.loc[:, params["frequency"], feature] +=_
→params["magnitude"] * np.linspace(0, 1, len(df_stress_test[:_
→params["frequency"]]))
        elif anomaly == "drop":
            df_stress_test.loc[:, params["frequency"], feature] =_
→params["magnitude"]
        elif anomaly == "noise":
            df_stress_test.loc[:, params["frequency"], feature] +=_
→params["magnitude"] * np.random.normal(0, 1, len(df_stress_test[:_
→params["frequency"]]))
```

Part 2: Evaluated model performance on stress-tested data using the aligned feature set.

```
[119]: from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
```

```

from sklearn.svm import SVC
from sklearn.naive_bayes import GaussianNB
from sklearn.tree import DecisionTreeClassifier
from xgboost import XGBClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score,
    f1_score

# Define models with best parameters
models = {
    "Random Forest": RandomForestClassifier(max_depth=10, max_features='sqrt',
    n_estimators=100, random_state=42),
    "Logistic Regression": LogisticRegression(C=0.1, solver='liblinear',
    random_state=42),
    "K-Nearest Neighbors": KNeighborsClassifier(n_neighbors=10, p=1,
    weights='distance'),
    "Support Vector Machine": SVC(C=0.1, kernel='rbf', random_state=42),
    "Gradient Boosting": GradientBoostingClassifier(learning_rate=0.05,
    max_depth=5, n_estimators=50, random_state=42),
    "XGBoost": XGBClassifier(learning_rate=0.01, max_depth=3, n_estimators=50,
    random_state=42),
    "Naive Bayes": GaussianNB(),
    "Decision Tree": DecisionTreeClassifier(criterion='gini', max_depth=10,
    min_samples_split=10, random_state=42)
}

# Results storage
stress_test_results = []

# Scale stress-tested data
X_stress_test_scaled = scaler.transform(df_stress_test[common_columns])

# Evaluate models on stressed data
for model_name, model in models.items():
    for anomaly_type in anomaly_types.keys():
        # Train on original data
        model.fit(X_original_scaled, y_original)

        # Predict on stress-tested data
        y_pred_stressed = model.predict(X_stress_test_scaled)

        # Calculate performance metrics
        accuracy = accuracy_score(df_stress_test['label'], y_pred_stressed)
        precision = precision_score(df_stress_test['label'], y_pred_stressed,
        zero_division=1)
        recall = recall_score(df_stress_test['label'], y_pred_stressed,
        zero_division=1)

```

```

f1 = f1_score(df_stress_test['label'], y_pred_stressed, zero_division=1)

# Append results
stress_test_results.append({
    "Model": model_name,
    "Anomaly Type": anomaly_type,
    "Accuracy": accuracy,
    "Precision": precision,
    "Recall": recall,
    "F1 Score": f1
})

# Convert results to DataFrame
df_stress_results = pd.DataFrame(stress_test_results)
print(df_stress_results)

```

	Model	Anomaly Type	Accuracy	Precision	Recall	F1 Score
0	Random Forest	spike	0.439983	0.007694	0.62	0.015200
1	Random Forest	drift	0.439983	0.007694	0.62	0.015200
2	Random Forest	drop	0.439983	0.007694	0.62	0.015200
3	Random Forest	noise	0.439983	0.007694	0.62	0.015200
4	Logistic Regression	spike	0.413286	0.007112	0.60	0.014056
5	Logistic Regression	drift	0.413286	0.007112	0.60	0.014056
6	Logistic Regression	drop	0.413286	0.007112	0.60	0.014056
7	Logistic Regression	noise	0.413286	0.007112	0.60	0.014056
8	K-Nearest Neighbors	spike	0.429667	0.007915	0.65	0.015640
9	K-Nearest Neighbors	drift	0.429667	0.007915	0.65	0.015640
10	K-Nearest Neighbors	drop	0.429667	0.007915	0.65	0.015640
11	K-Nearest Neighbors	noise	0.429667	0.007915	0.65	0.015640
12	Support Vector Machine	spike	0.346647	0.010556	1.00	0.020892
13	Support Vector Machine	drift	0.346647	0.010556	1.00	0.020892
14	Support Vector Machine	drop	0.346647	0.010556	1.00	0.020892
15	Support Vector Machine	noise	0.346647	0.010556	1.00	0.020892
16	Gradient Boosting	spike	0.448209	0.008056	0.64	0.015912
17	Gradient Boosting	drift	0.448209	0.008056	0.64	0.015912
18	Gradient Boosting	drop	0.448209	0.008056	0.64	0.015912
19	Gradient Boosting	noise	0.448209	0.008056	0.64	0.015912
20	XGBoost	spike	0.375087	0.007555	0.68	0.014943
21	XGBoost	drift	0.375087	0.007555	0.68	0.014943
22	XGBoost	drop	0.375087	0.007555	0.68	0.014943
23	XGBoost	noise	0.375087	0.007555	0.68	0.014943
24	Naive Bayes	spike	0.637669	0.018512	0.98	0.036337
25	Naive Bayes	drift	0.637669	0.018512	0.98	0.036337
26	Naive Bayes	drop	0.637669	0.018512	0.98	0.036337
27	Naive Bayes	noise	0.637669	0.018512	0.98	0.036337
28	Decision Tree	spike	0.419350	0.008246	0.69	0.016297
29	Decision Tree	drift	0.419350	0.008246	0.69	0.016297

30	Decision Tree	drop	0.419350	0.008246	0.69	0.016297
31	Decision Tree	noise	0.419350	0.008246	0.69	0.016297

Part 3: Visualized stress testing results by model and anomaly type.

```
[120]: import seaborn as sns
import matplotlib.pyplot as plt

# Plot metrics for stress testing results
fig, ax = plt.subplots(2, 2, figsize=(16, 12))
metrics = ["Accuracy", "Precision", "Recall", "F1 Score"]
axes = ax.ravel()

for i, metric in enumerate(metrics):
    sns.barplot(data=df_stress_results, x="Model", y=metric, hue="Anomaly Type", ax=axes[i])
    axes[i].set_title(f"Model Performance on Stress Test - {metric}")
    axes[i].set_xticklabels(axes[i].get_xticklabels(), rotation=45, ha='right')

plt.tight_layout()
plt.show()
```

