

1. Introduction

This report details the process of anomaly detection in the provided dataset using the Isolation Forest algorithm. The dataset, which doesn't have a specific name and contains 61 features (x1 to x60, plus an unnamed column) and a binary response variable (y), has been analyzed for anomalies. The primary objective was to identify unusual data points that deviate significantly from the majority of the data.

2. Data Exploration and Preprocessing

Initial exploration of the data involved examining the first few rows and the data types of each column. The dataset was checked for missing values, but none were found. Subsequently, a descriptive analysis was performed, providing statistical summaries like mean, standard deviation, minimum, and maximum for each column. Additionally, the number of unique values for each column was determined.

	y	x1	x2	x3	x4	x5
count	18398	18398	18398	18398	18398	18398
mean	0.00673986	0.0118235	0.157986	0.5693	-9.95835	0.00673986
std	0.0818218	0.742875	4.93976	5.93718	131.034	0.634018
min	0	-3.78728	-17.3165	-18.1985	-322.782	-1.62395
25%	0	-0.405681	-2.15823	-3.53705	-111.378	-0.446018
50%	0	0.128245	-0.0755048	-0.190683	-14.8816	-0.120018

Number of duplicate rows: 0

Descriptive Statistics:

	time	y	x1 \
count	18398	18398.000000	18398.000000
mean	1999-05-15 01:20:42.728557312	0.006740	0.011824
min	1999-05-01 00:00:00	0.000000	-3.787279
25%	1999-05-08 03:36:30	0.000000	-0.405681
50%	1999-05-14 18:39:00	0.000000	0.128245
75%	1999-05-22 06:01:30	0.000000	0.421222
max	1999-05-29 00:06:00	1.000000	3.054156

std	NaN	0.081822	0.742875
-----	-----	----------	----------

	x2	x3	x4	x5	x6 \
count	18398.000000	18398.000000	18398.000000	18398.000000	18398.000000
mean	0.157986	0.569300	-9.958345	0.006518	2.387533
min	-17.316550	-18.198509	-322.781610	-1.623988	-279.408440
25%	-2.158235	-3.537054	-111.378372	-0.446787	-24.345268
50%	-0.075505	-0.190683	-14.881585	-0.120745	10.528435
75%	2.319297	3.421223	92.199134	0.325152	32.172974
max	16.742105	15.900116	334.694098	4.239385	96.060768
std	4.939762	5.937178	131.033712	0.634054	37.104012

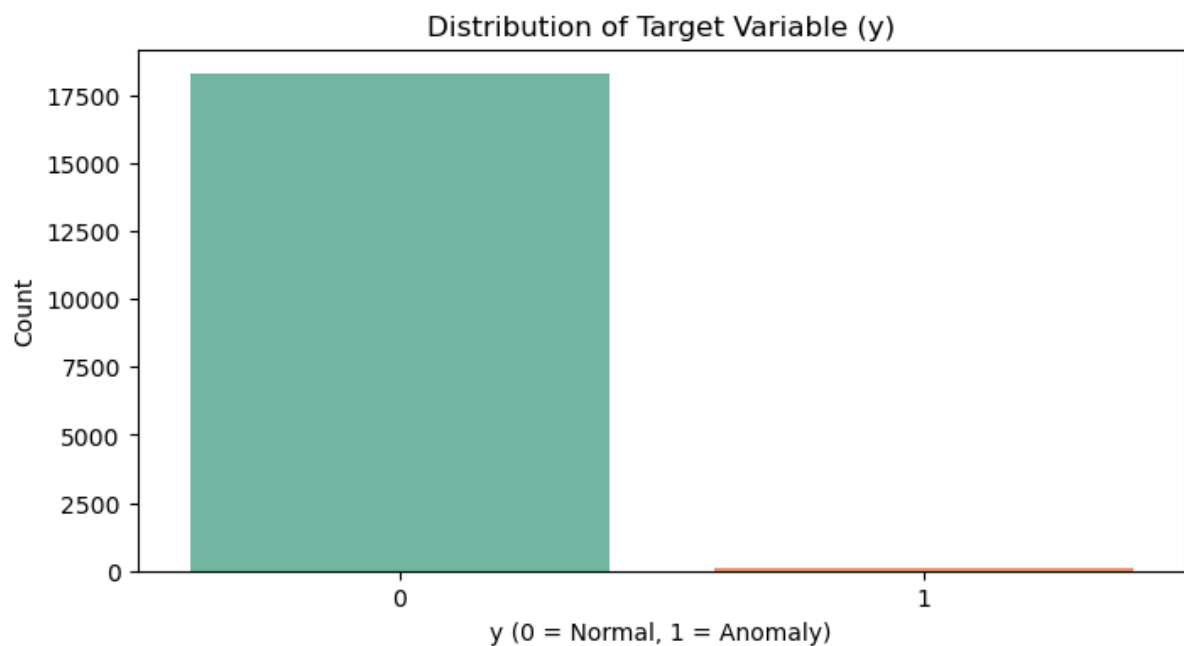
	x7	x8 ...	x51	x52 \
count	18398.000000	18398.000000	...	18398.000000 18398.000000
mean	0.001647	-0.004125	...	-3.357339 0.380519
min	-0.429273	-0.451141	...	-3652.989000 -187.943440
25%	-0.058520	-0.051043	...	29.984624 -3.672684
50%	-0.009338	-0.000993	...	29.984624 0.294846
75%	0.060515	0.038986	...	29.984624 5.109543
max	1.705590	0.788826	...	40.152348 14.180588
std	0.108870	0.075460	...	348.256716 6.211598

	x54	x55	x56	x57	x58 \
count	18398.000000	18398.000000	18398.000000	18398.000000	18398.000000
mean	0.173708	2.379154	9.234953	0.233493	-0.001861
min	-8.210370	-230.574030	-269.039500	-12.640370	-0.149790
25%	0.487780	-40.050046	-45.519149	-1.598804	0.000470
50%	0.702299	17.471317	1.438806	0.085826	0.012888
75%	2.675751	44.093387	63.209681	2.222118	0.020991

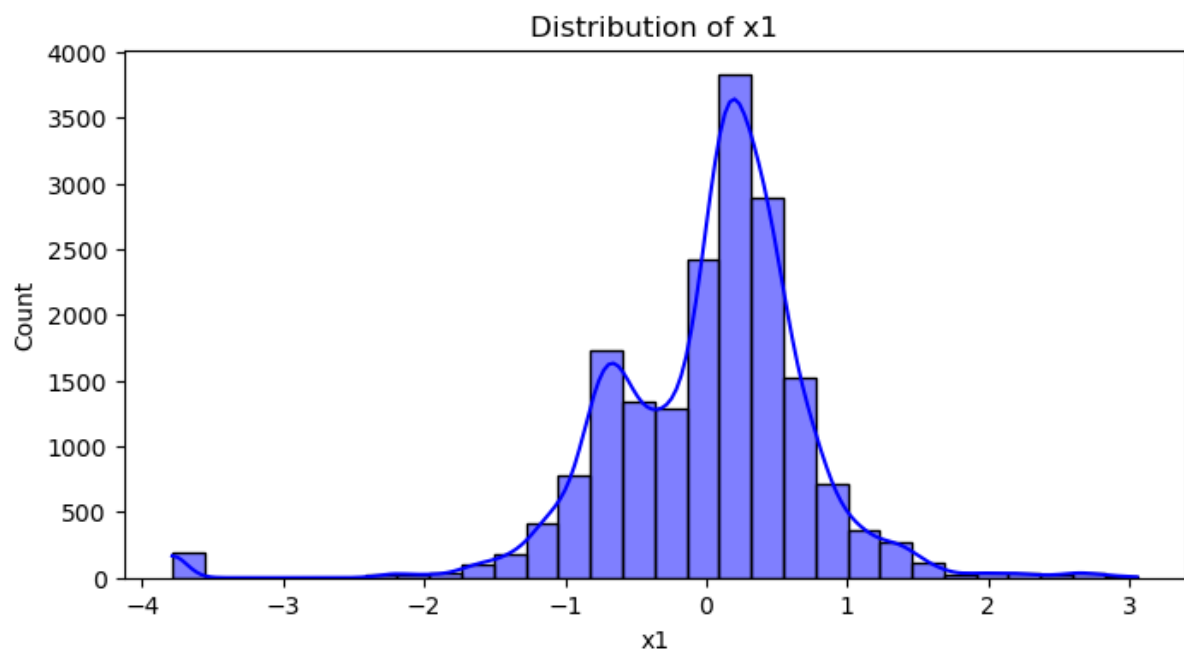
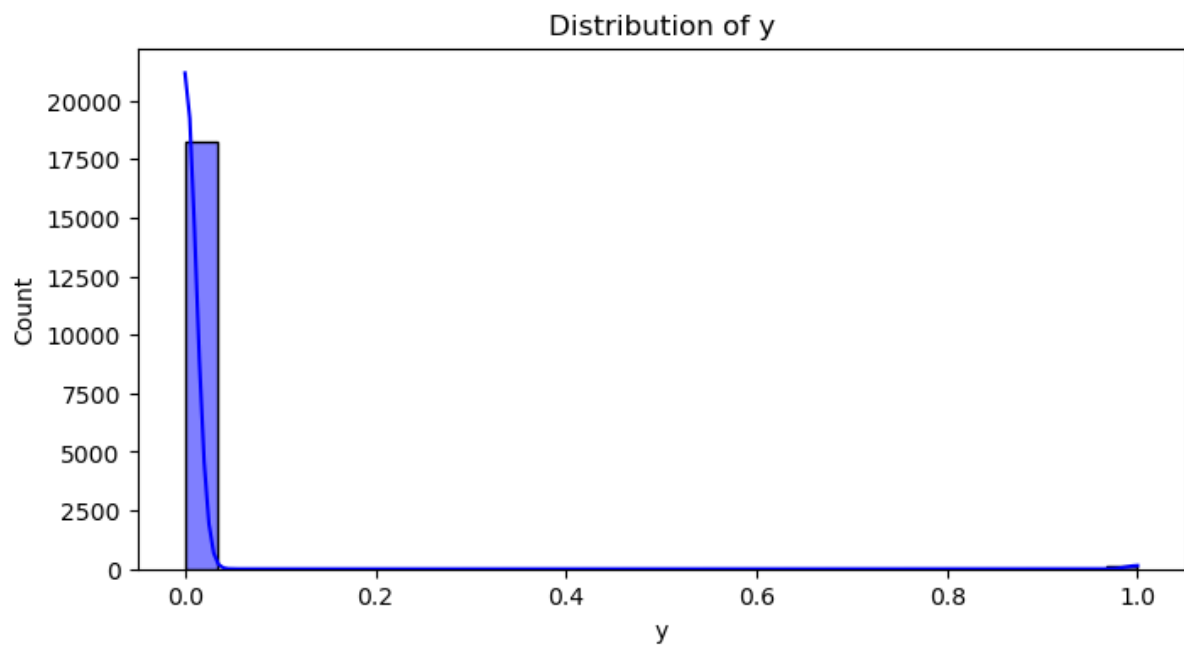
max	6.637265	287.252017	252.147455	6.922008	0.067249
std	3.029516	67.940694	81.274103	2.326838	0.048732

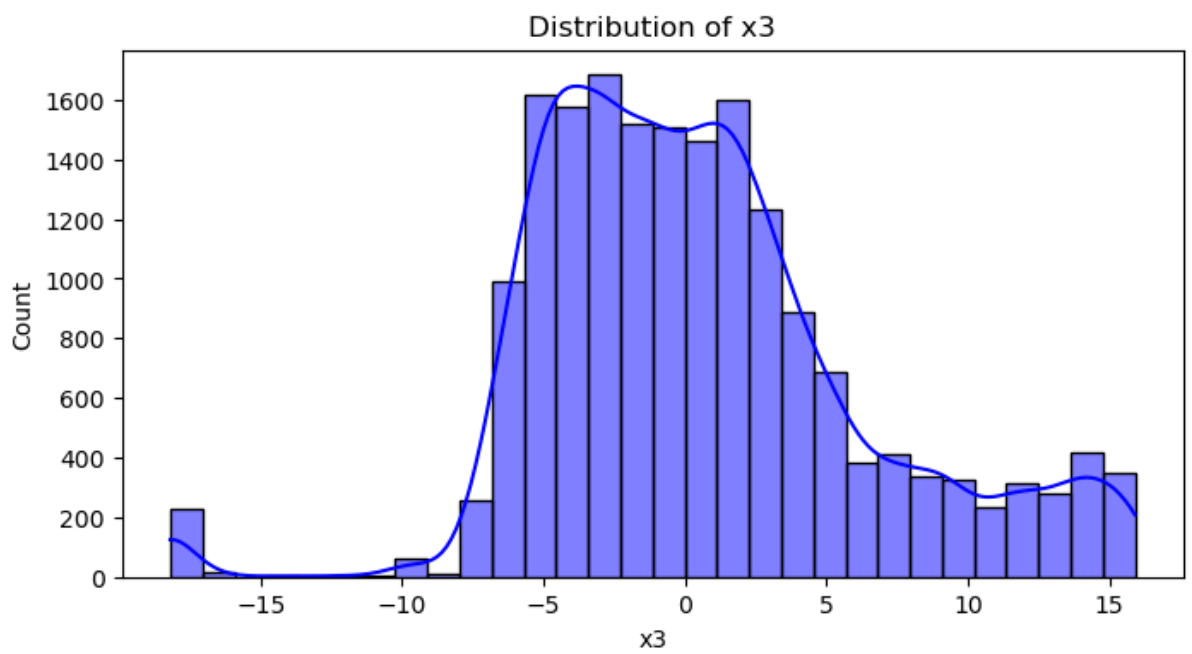
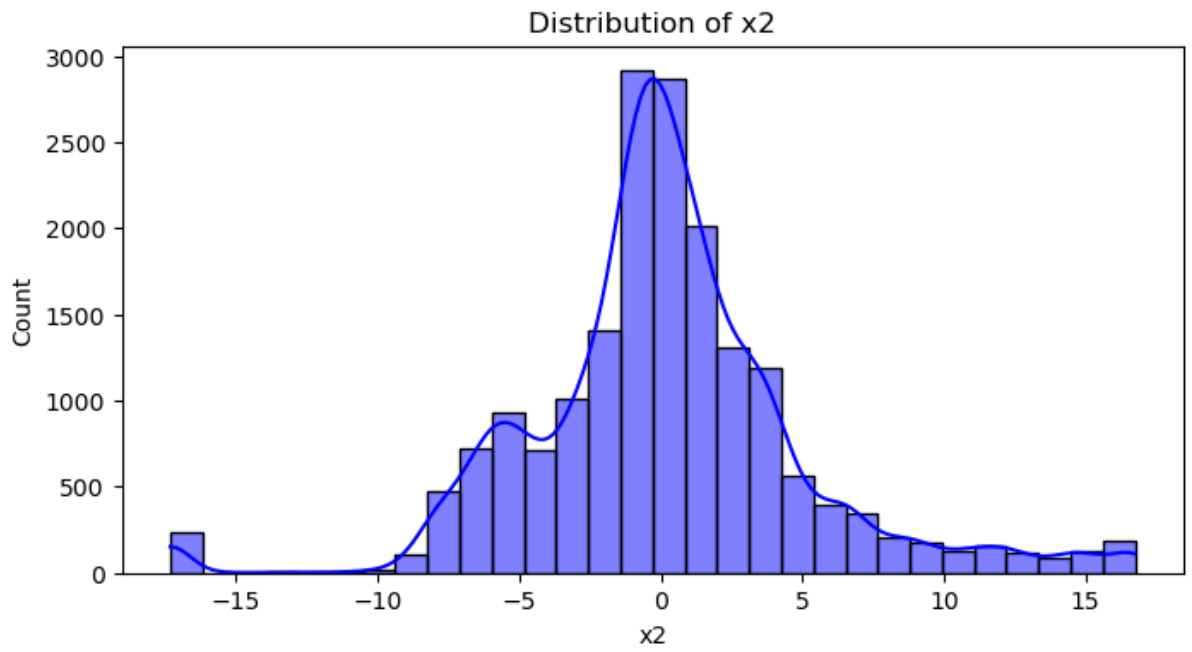
	x59	x60	y.1
count	18398.000000	18398.000000	18398.000000
mean	-0.061522	0.001258	0.001033
min	-100.810500	-0.012229	0.000000
25%	0.295023	-0.001805	0.000000
50%	0.734591	0.000710	0.000000
75%	1.266506	0.004087	0.000000
max	6.985460	0.020510	1.000000
std	10.394085	0.004721	0.032120

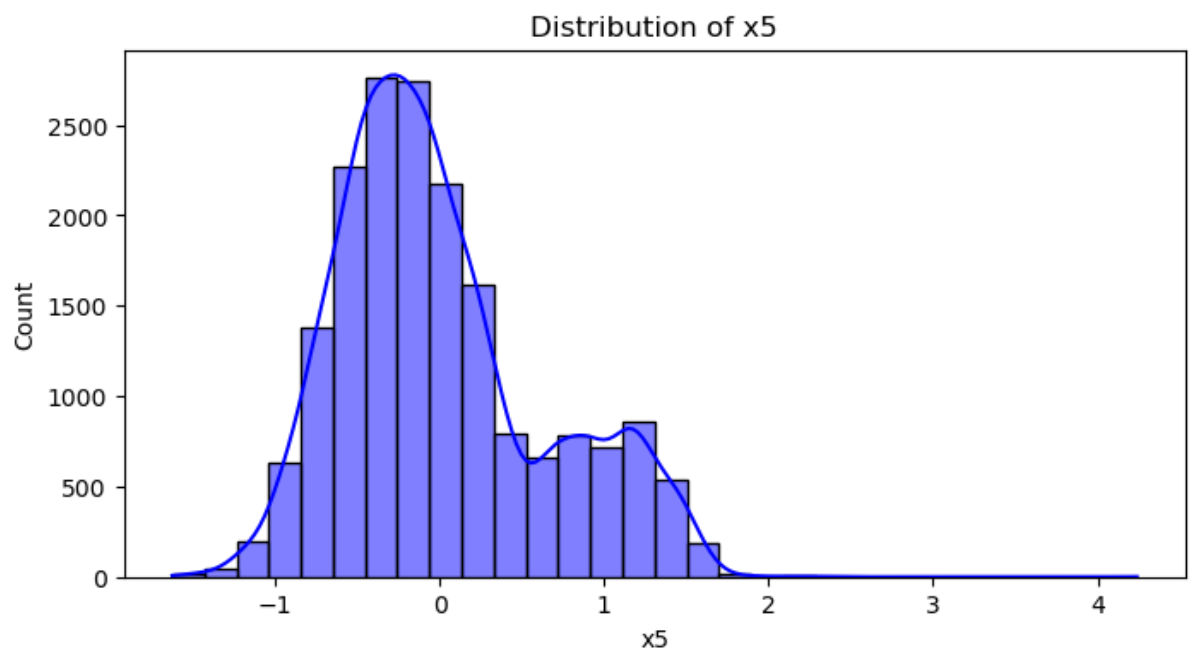
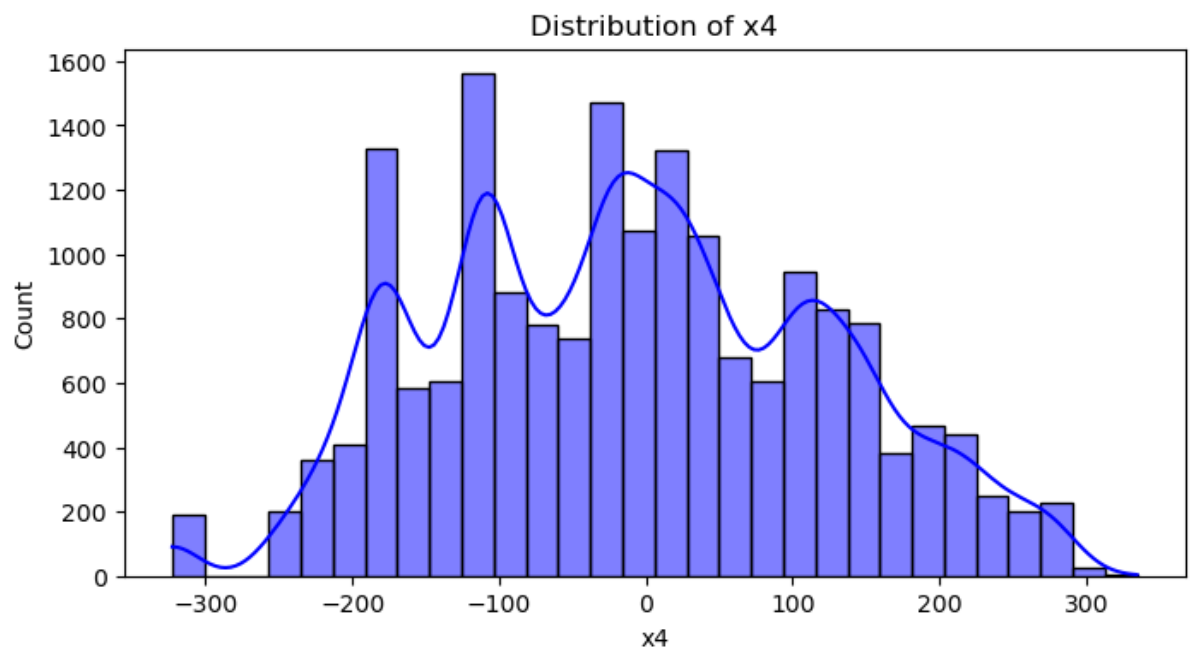
[8 rows x 62 columns]

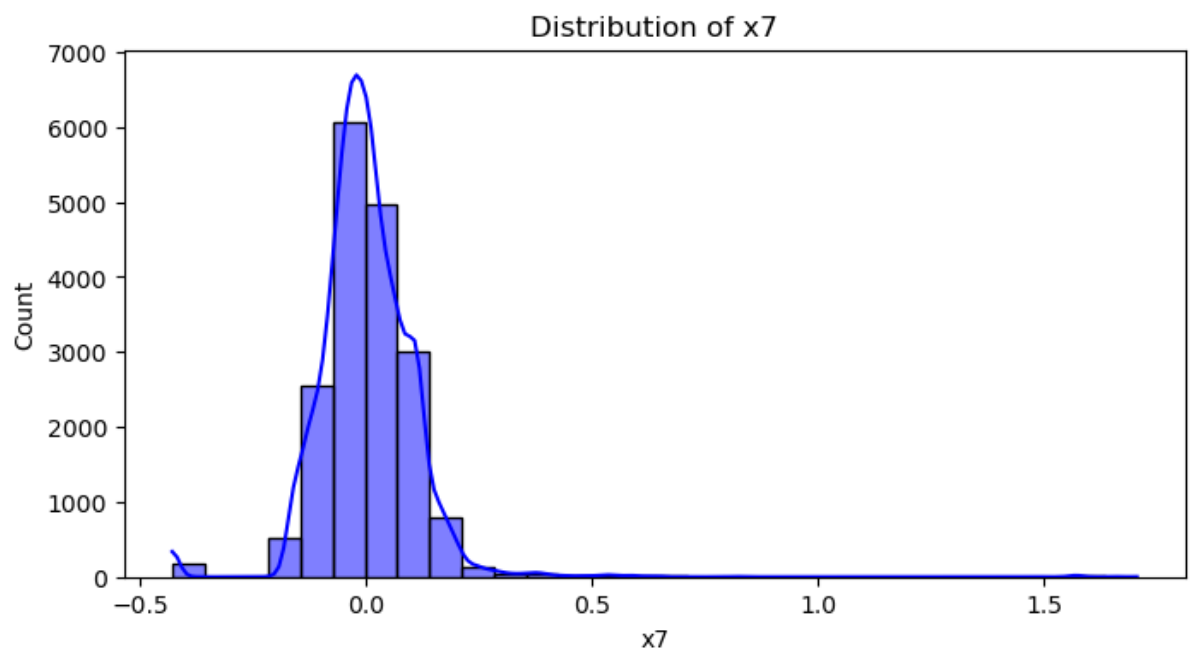
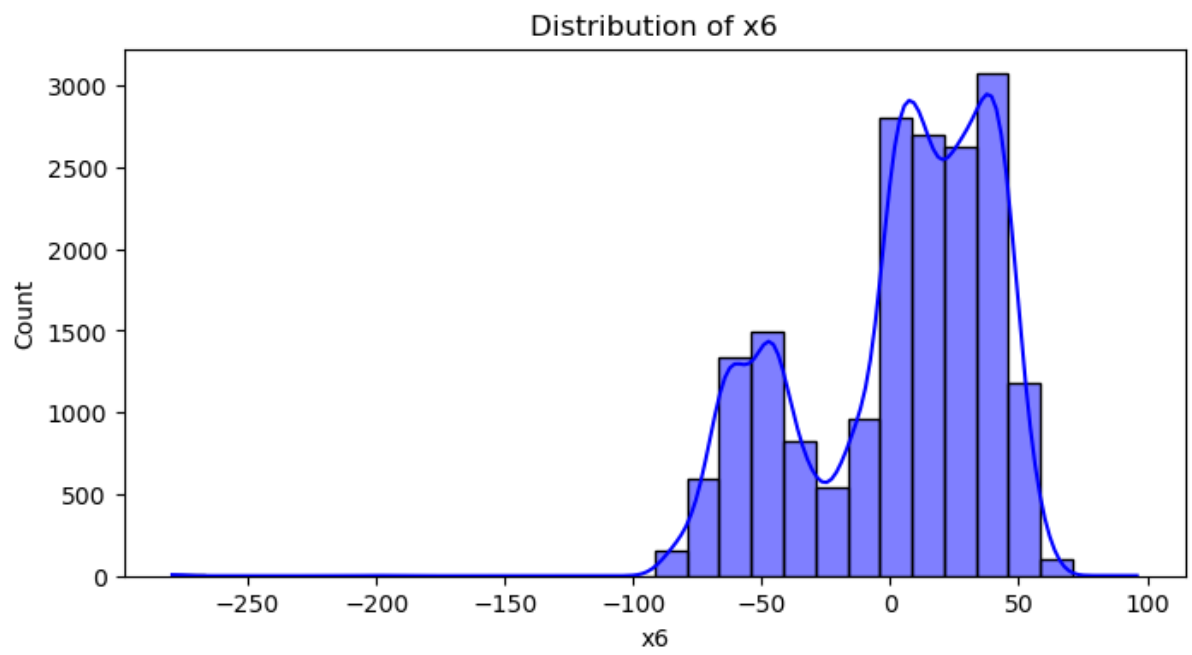


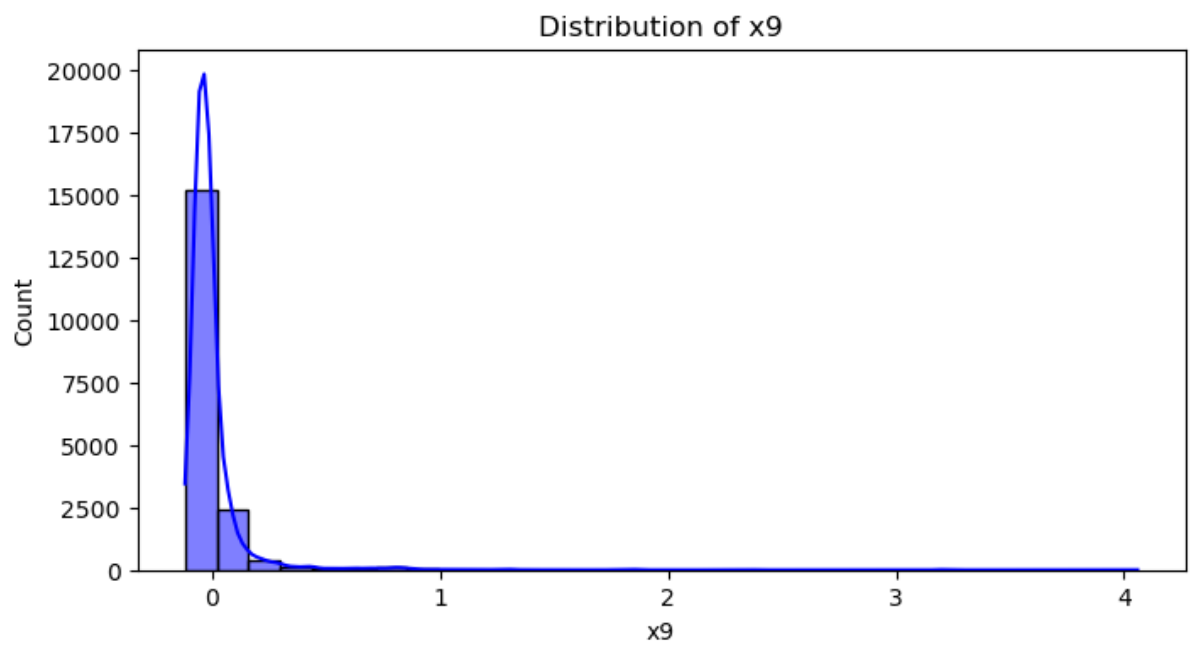
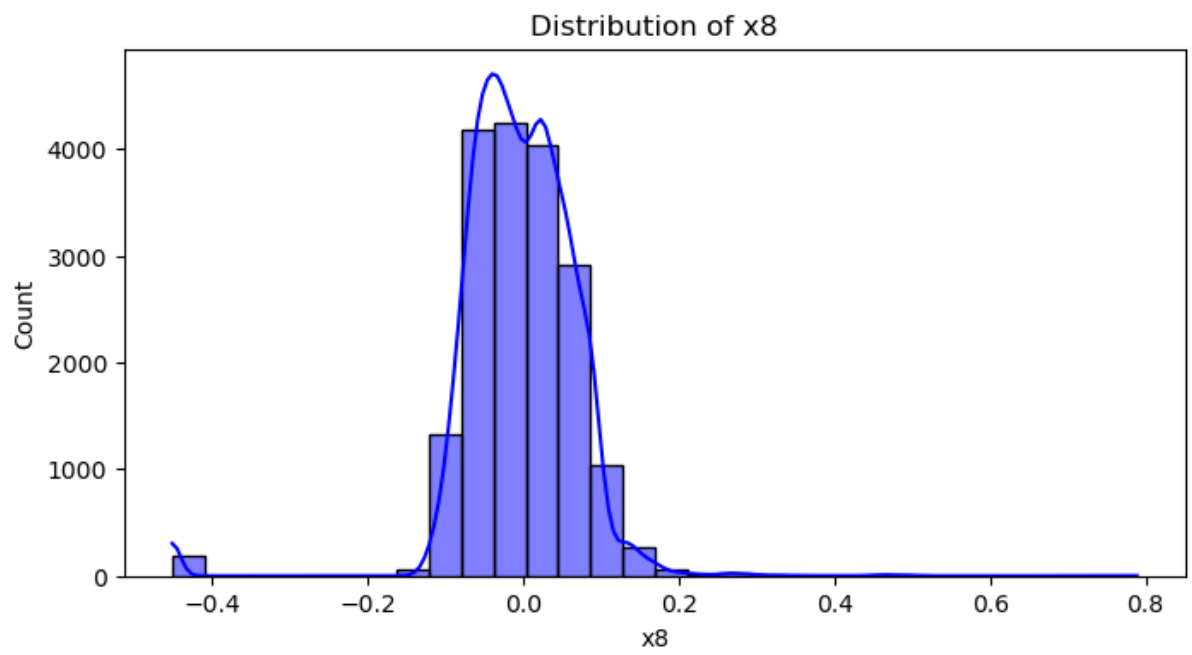
Generating histograms for numeric features:



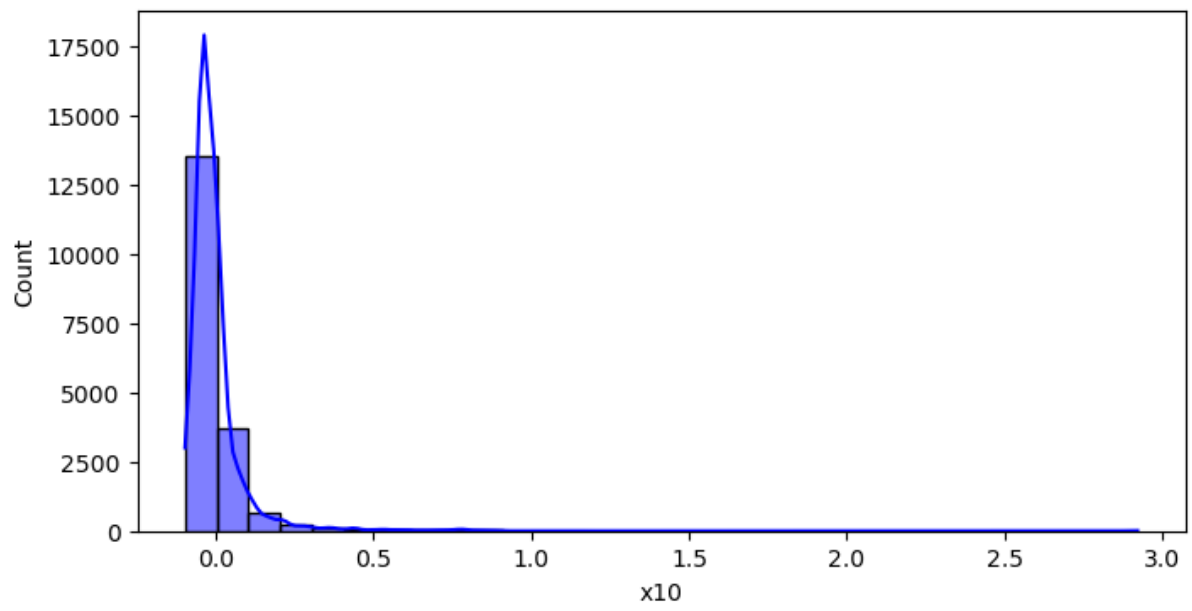




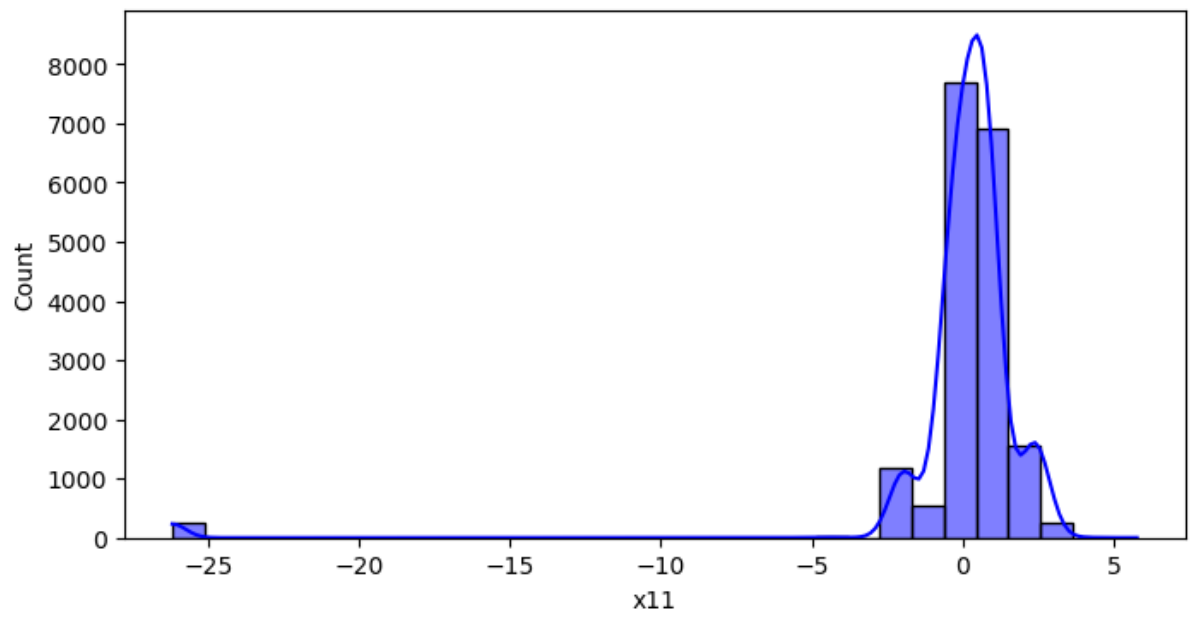


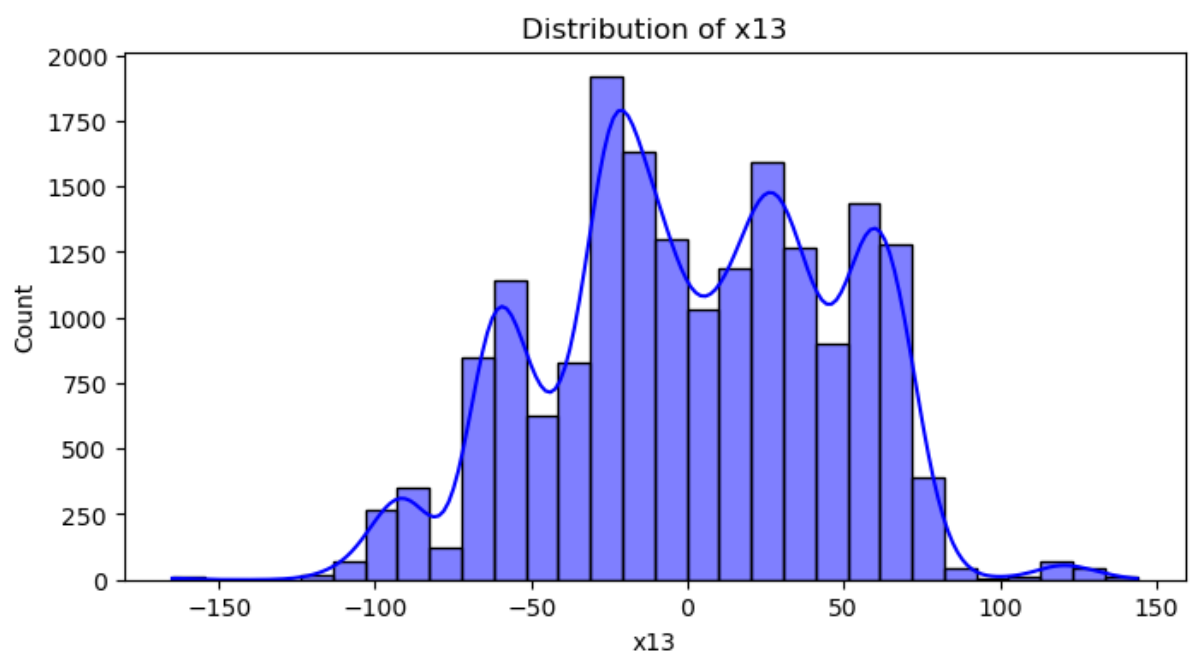
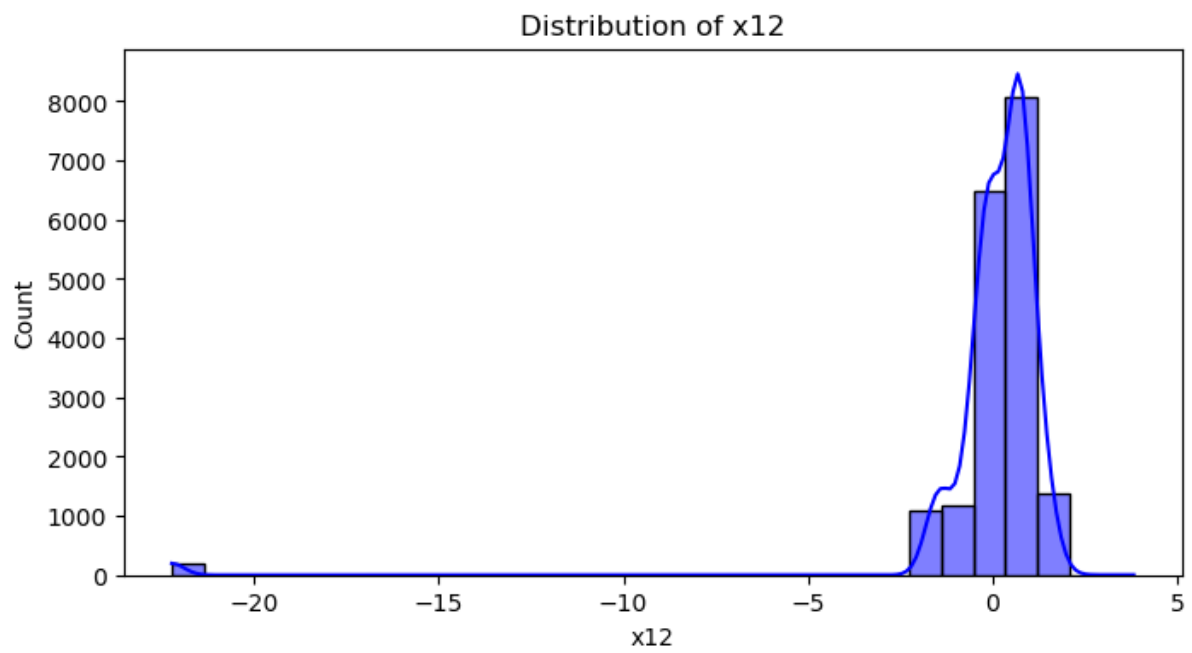


Distribution of x10

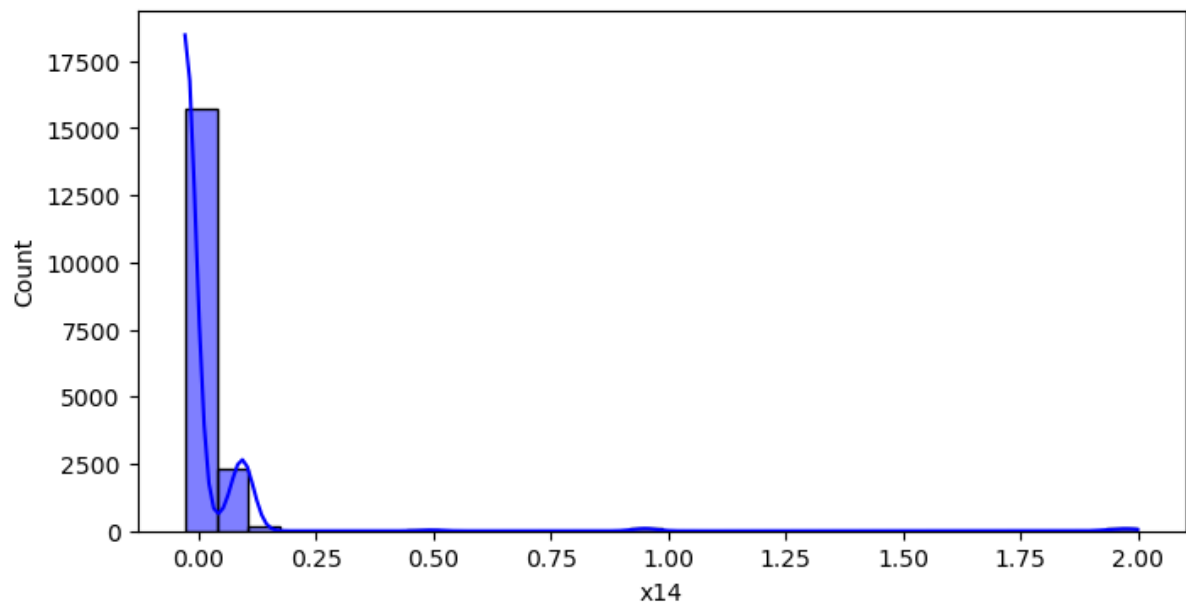


Distribution of x11

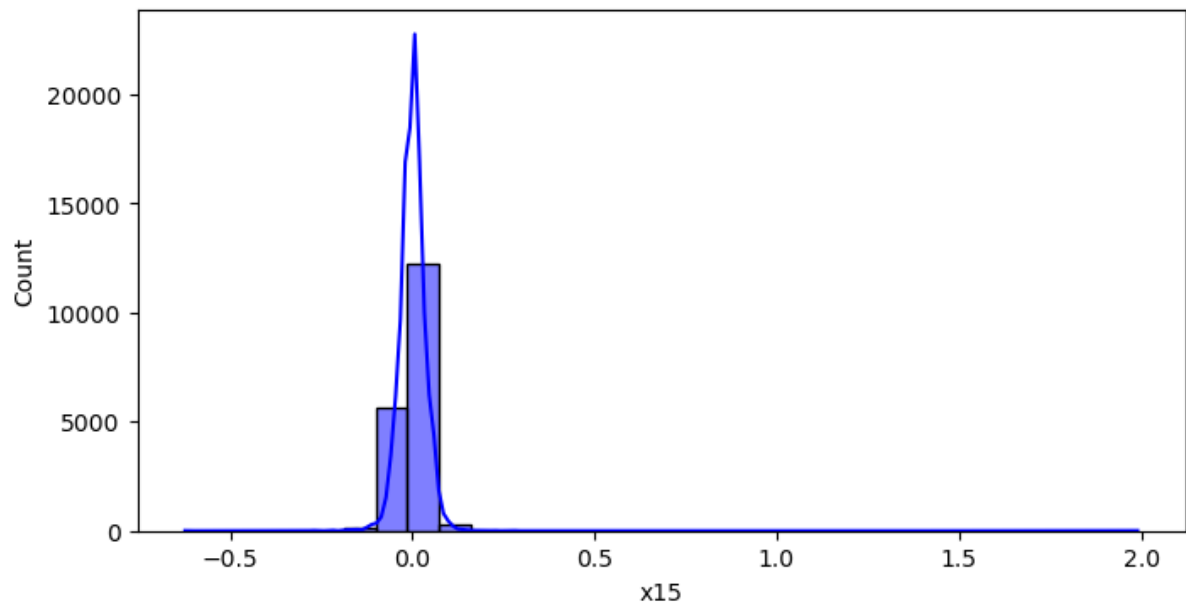


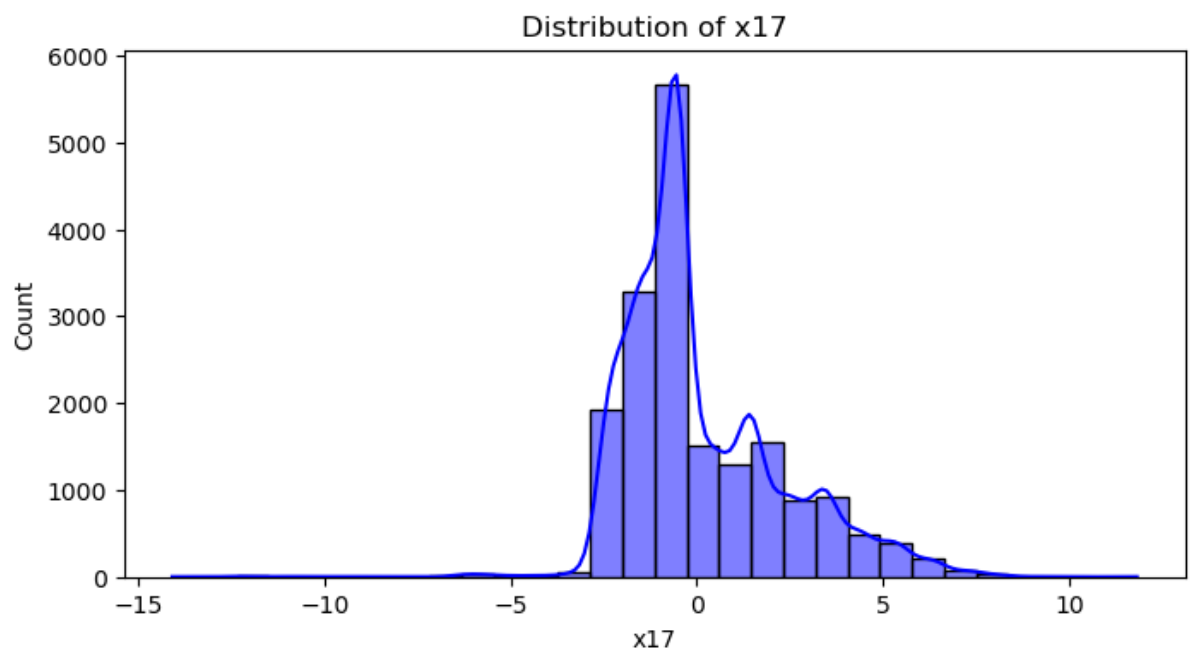
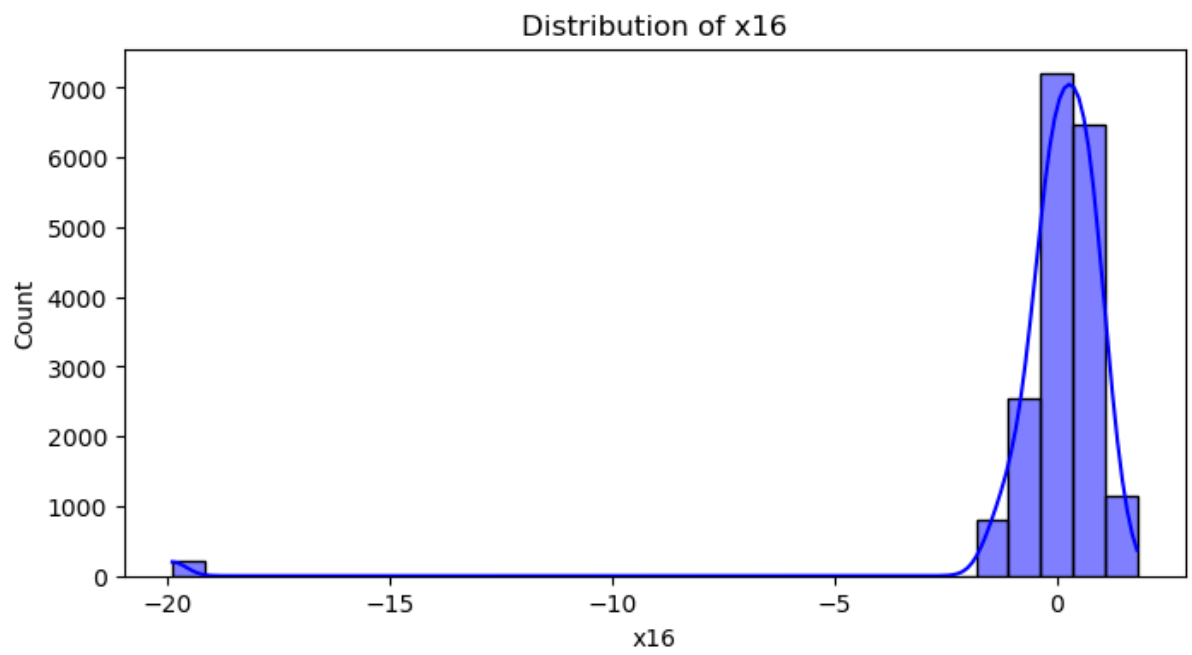


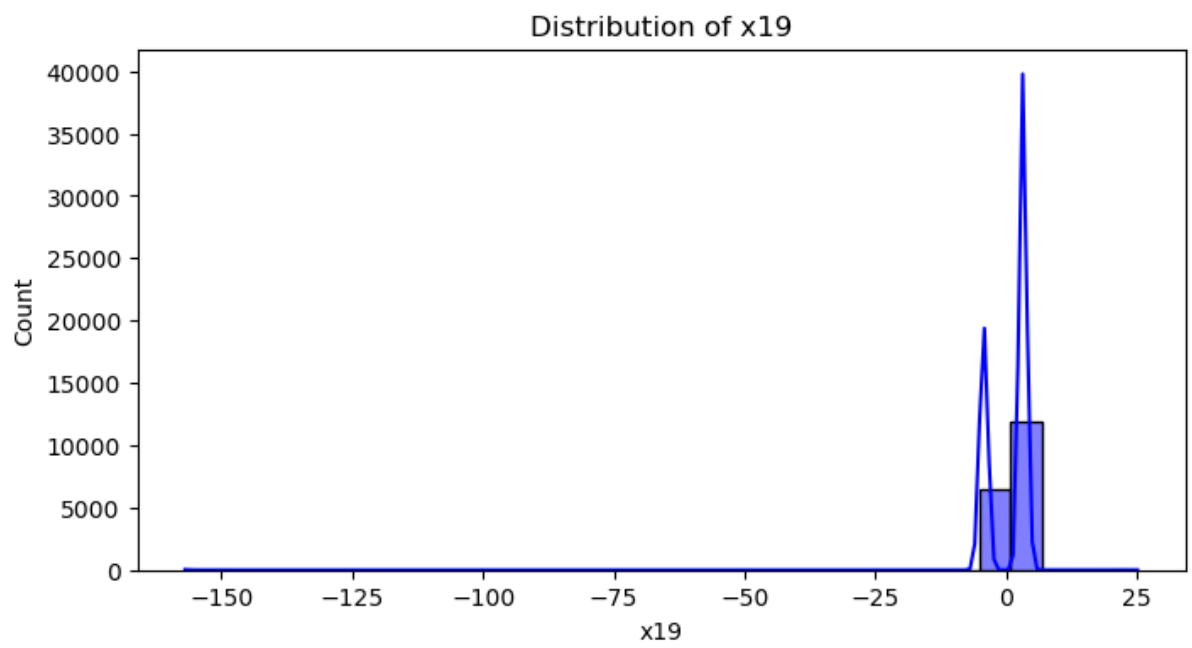
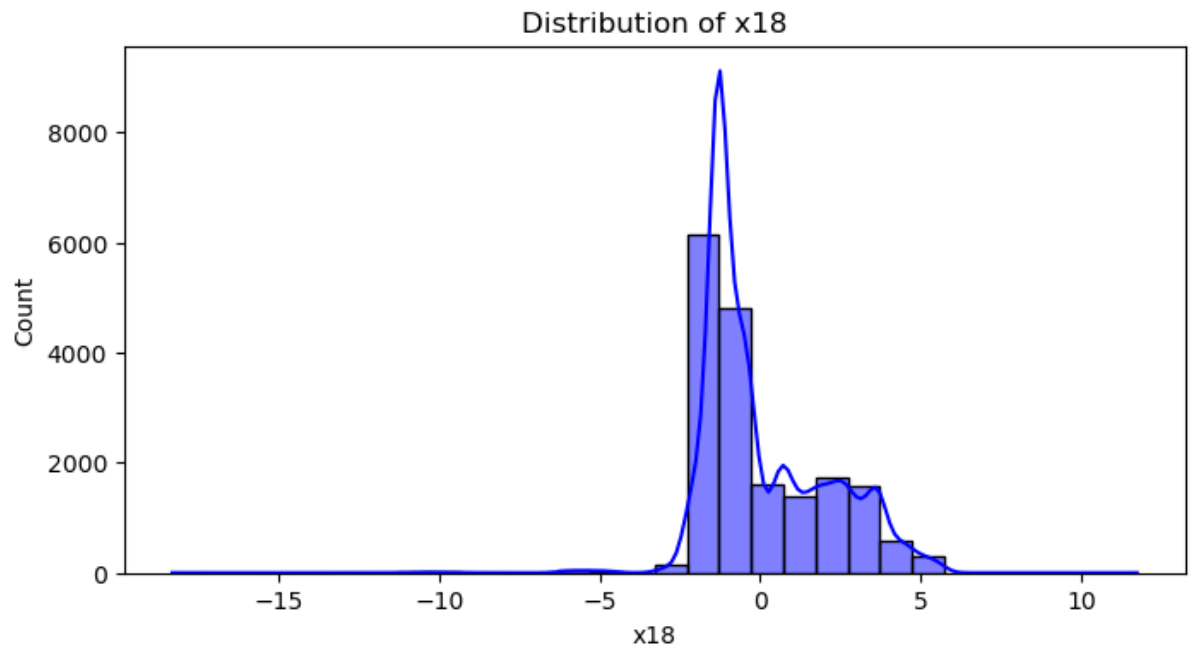
Distribution of x14



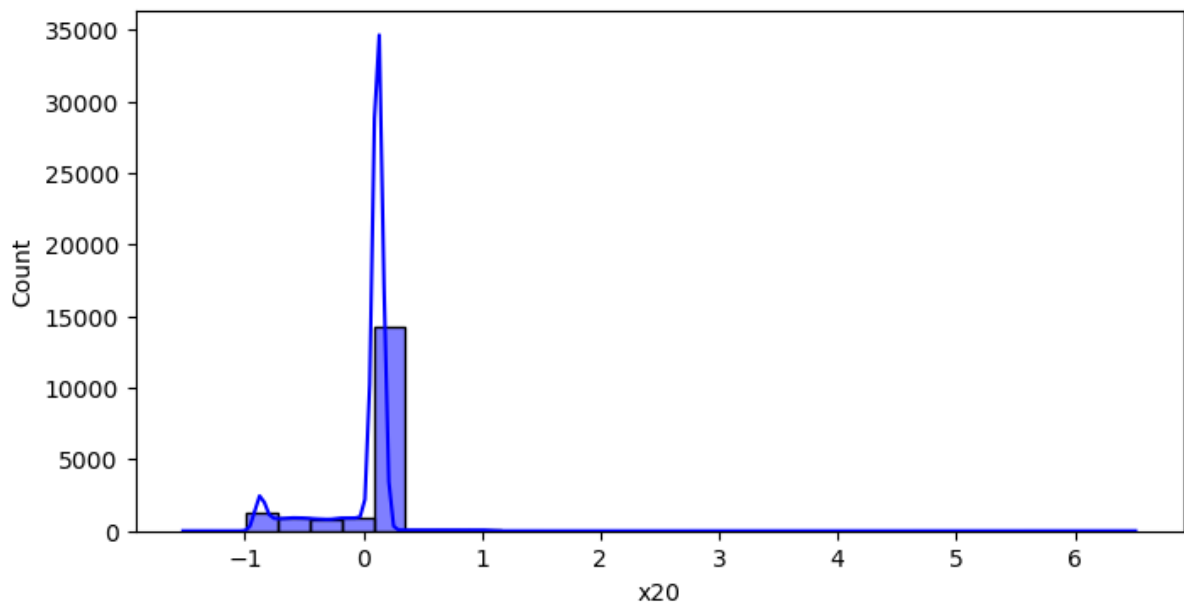
Distribution of x15



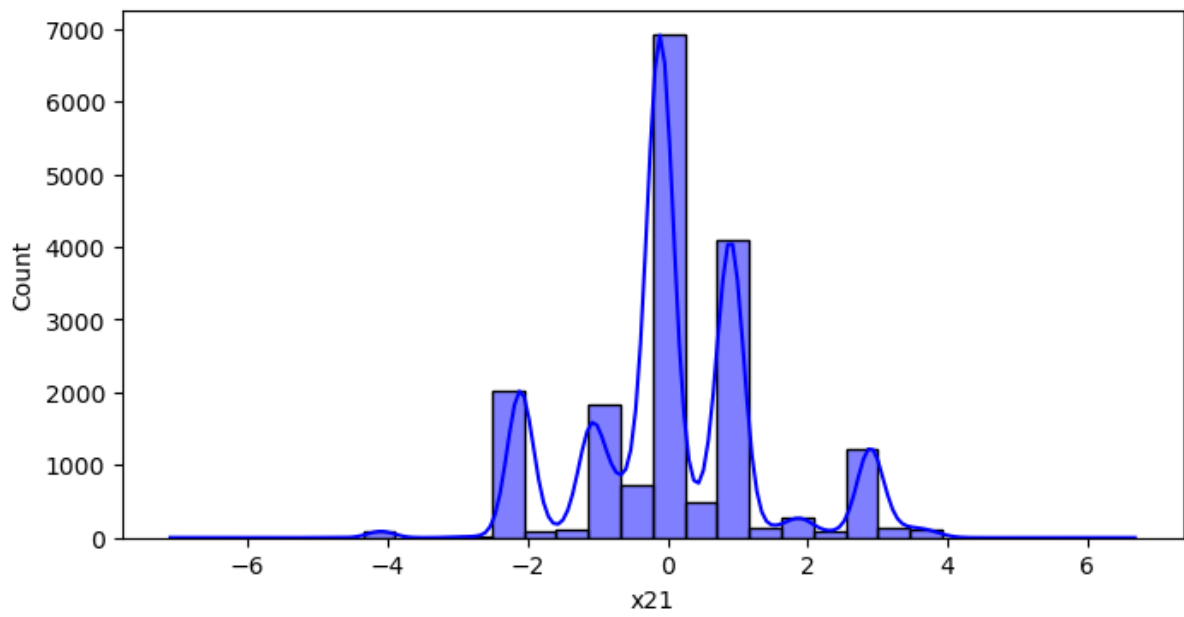


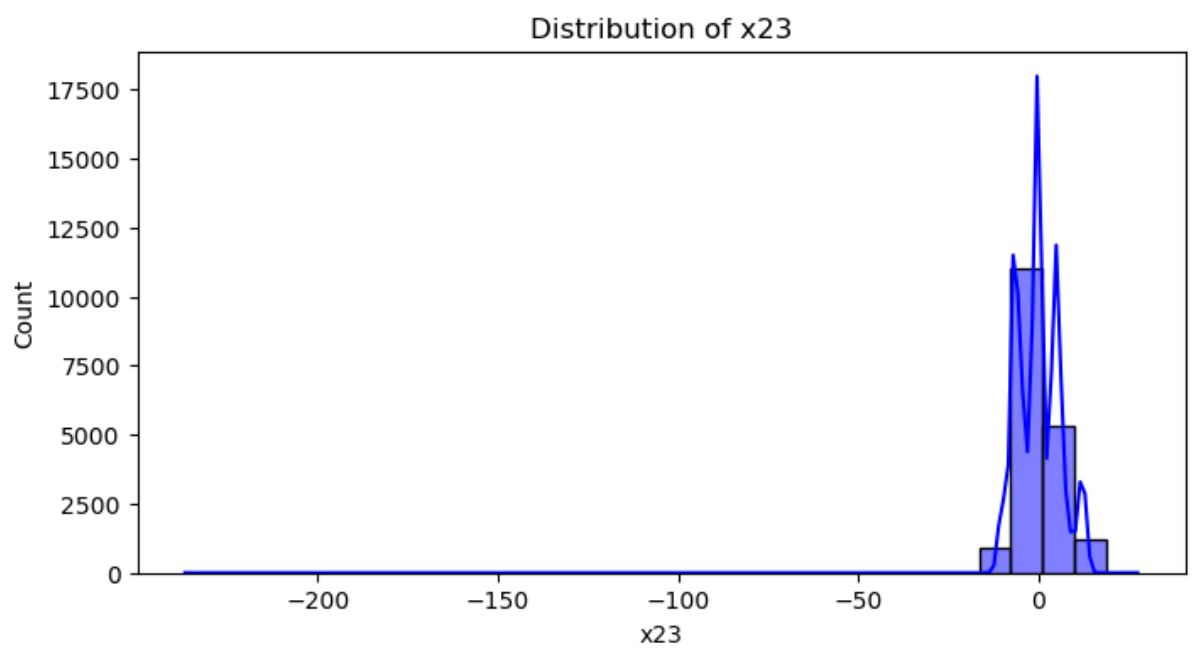
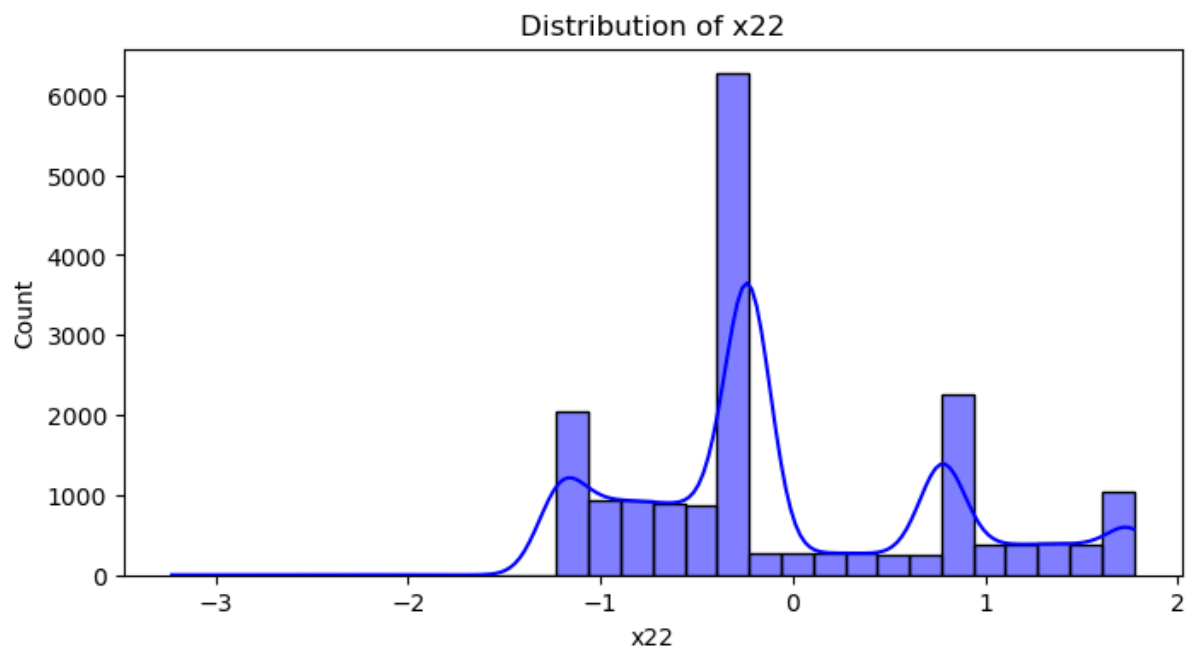


Distribution of x20

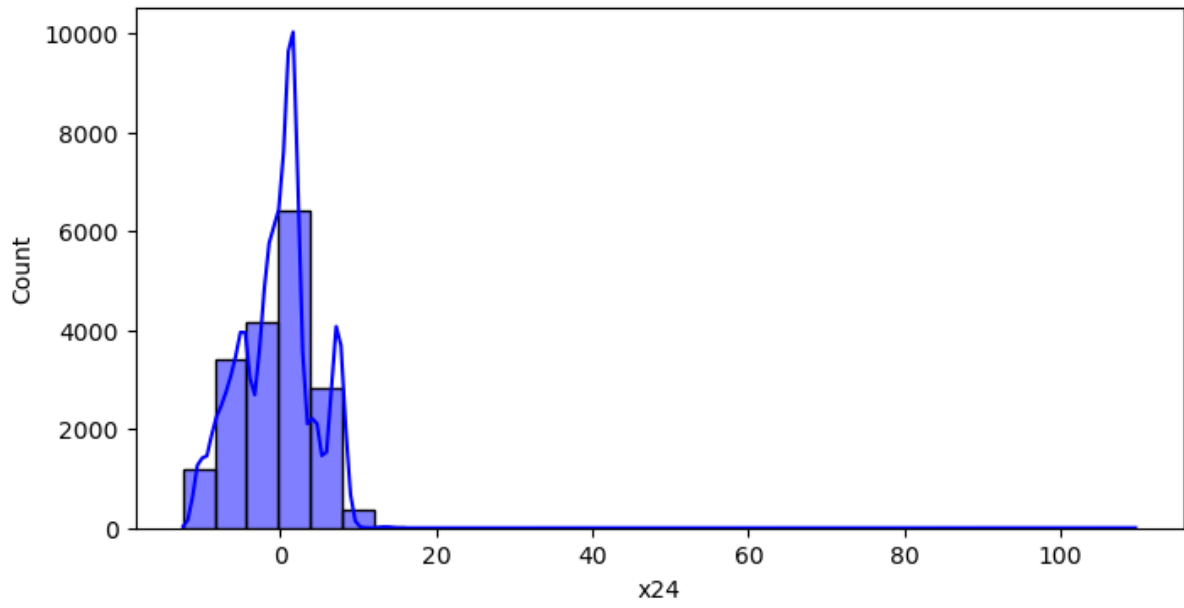


Distribution of x21

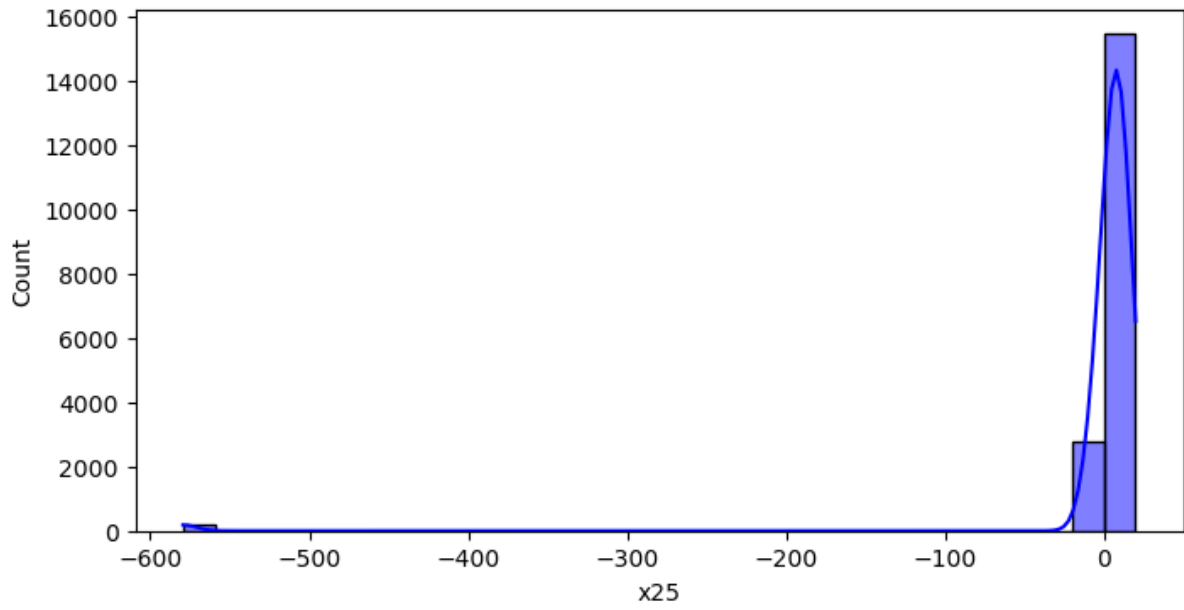




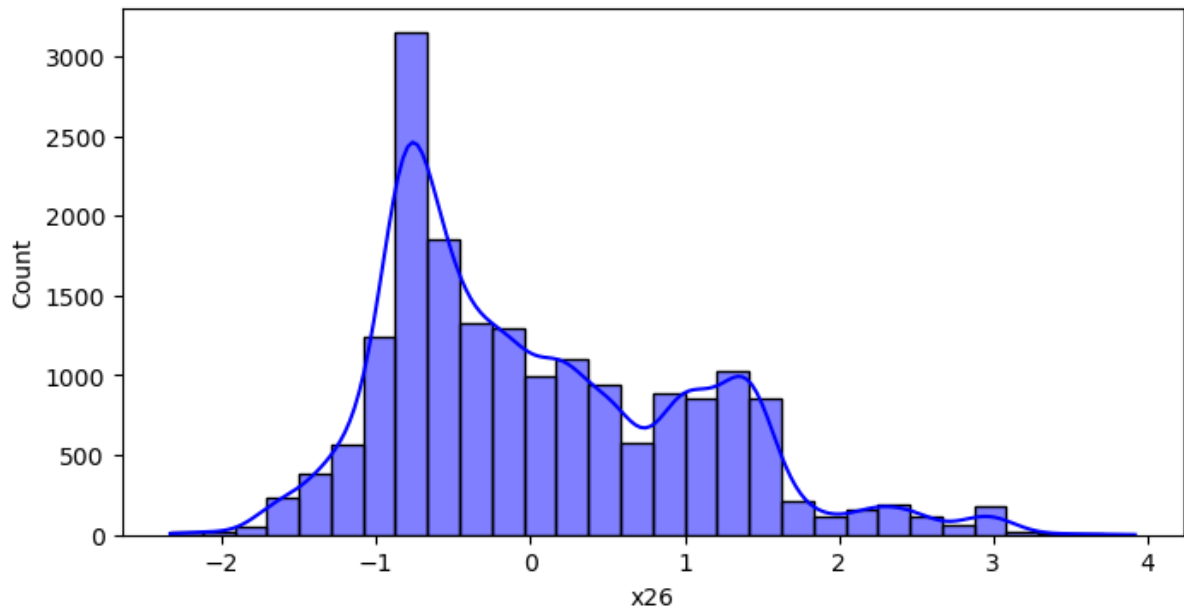
Distribution of x24



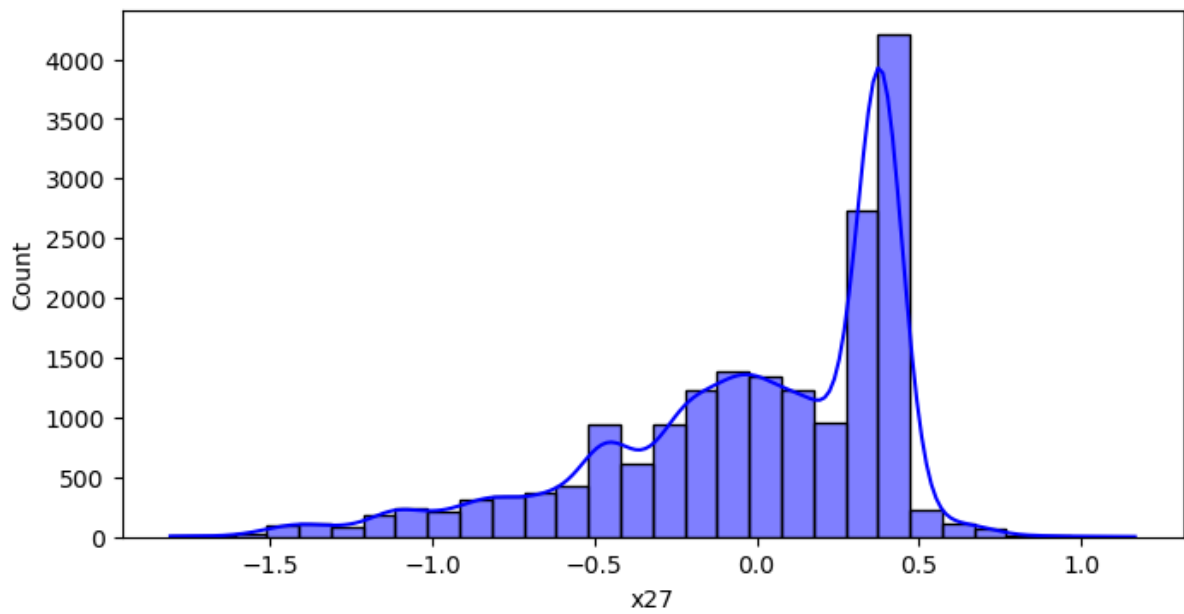
Distribution of x25

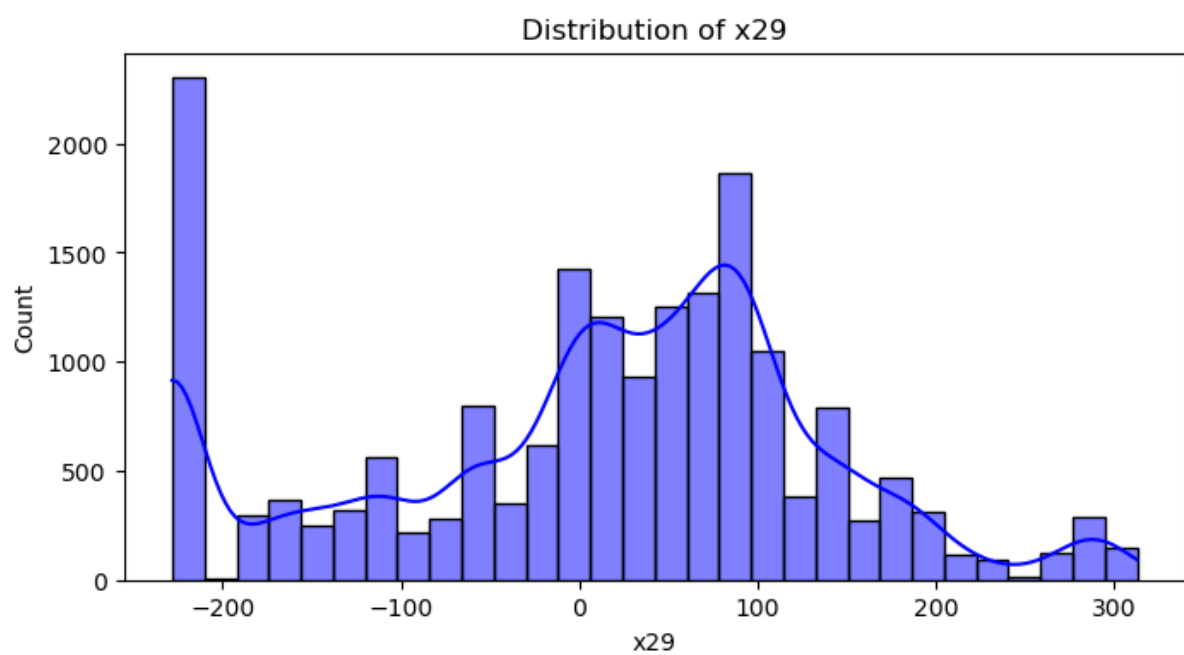
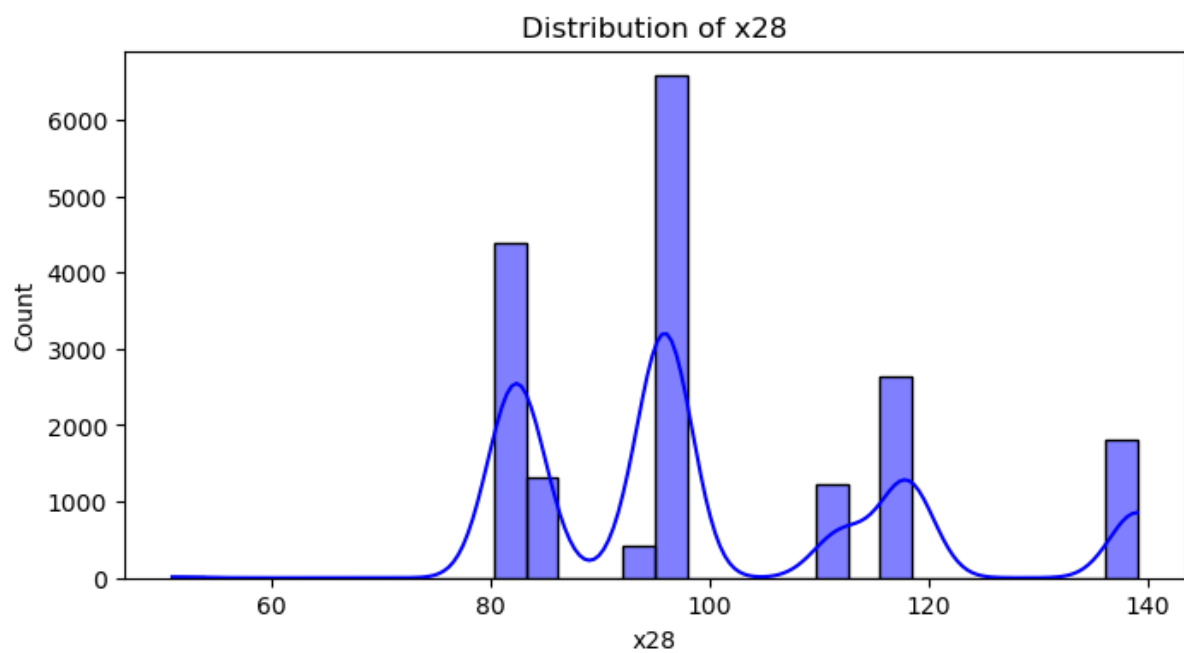


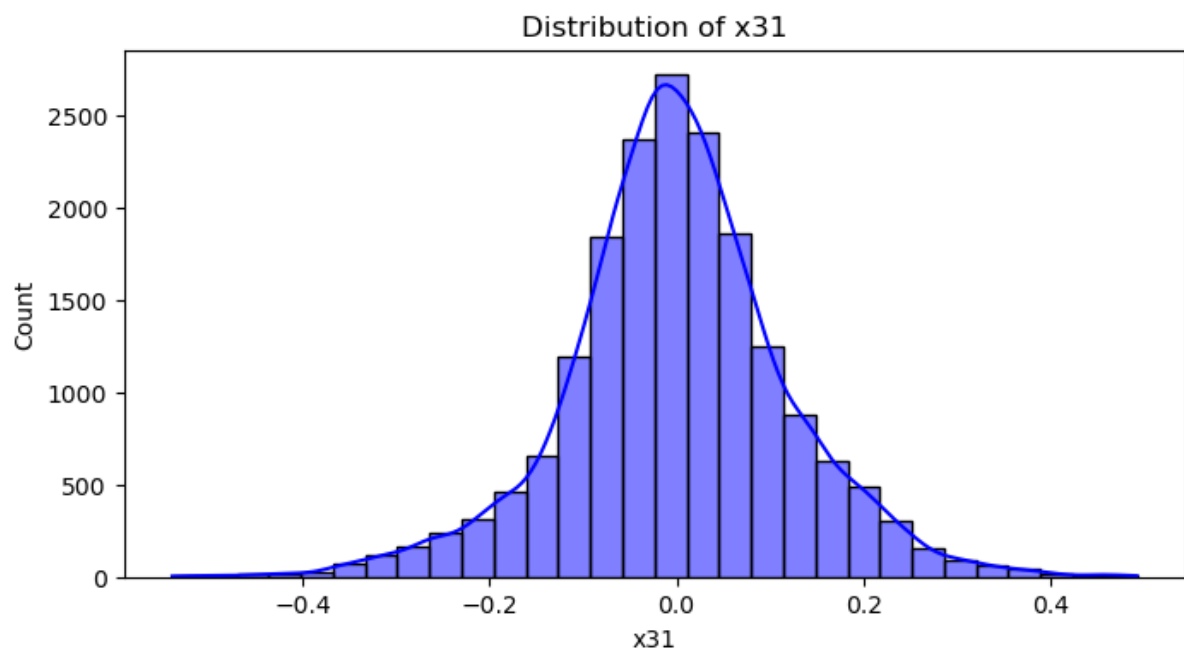
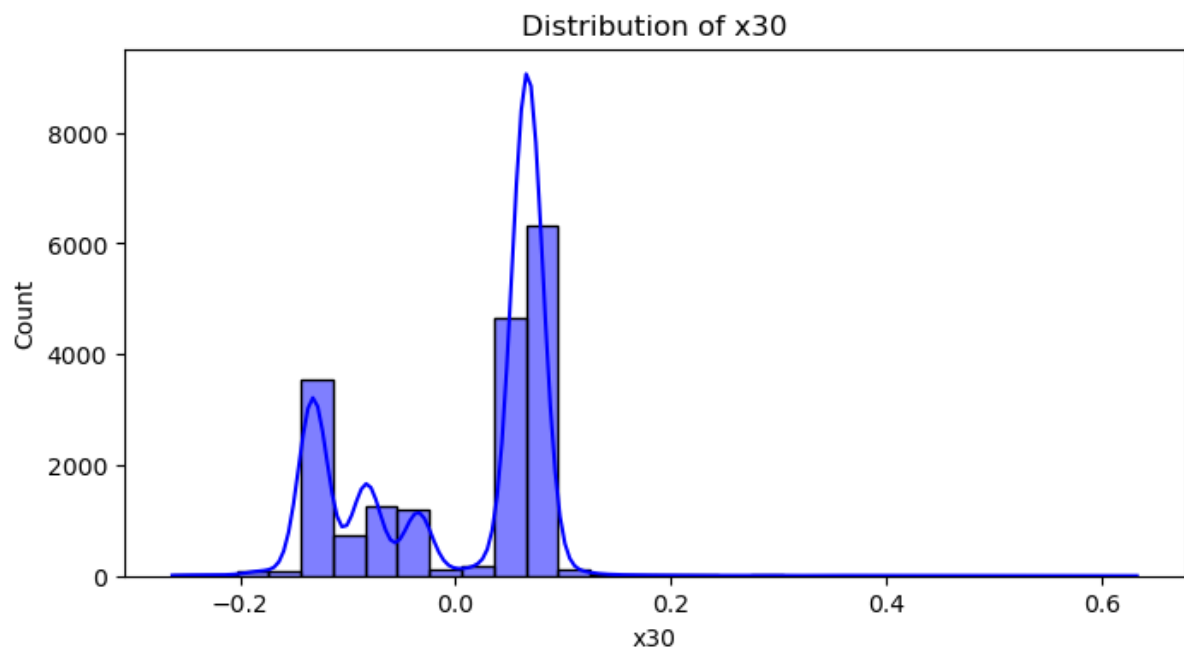
Distribution of x26

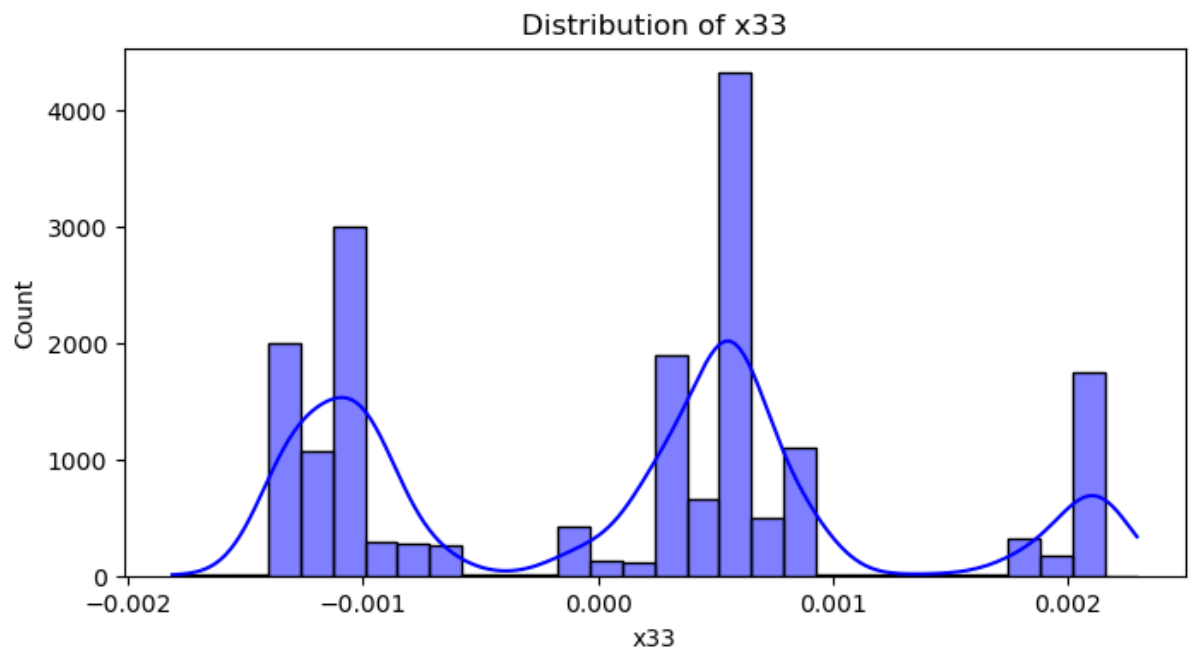
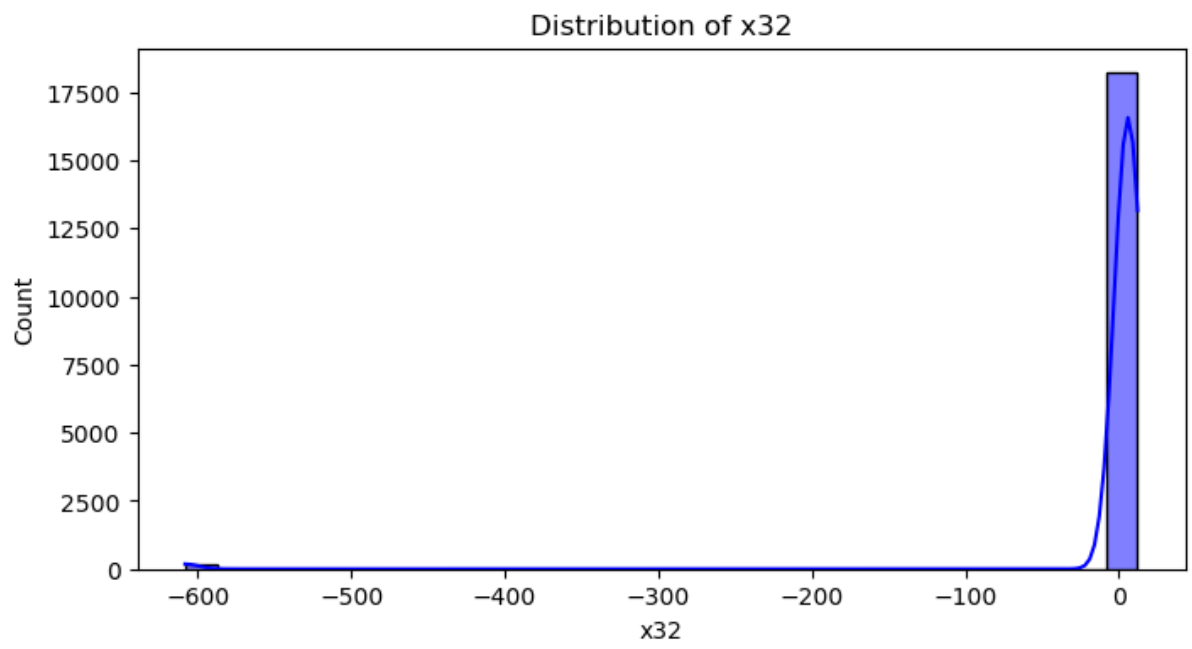


Distribution of x27

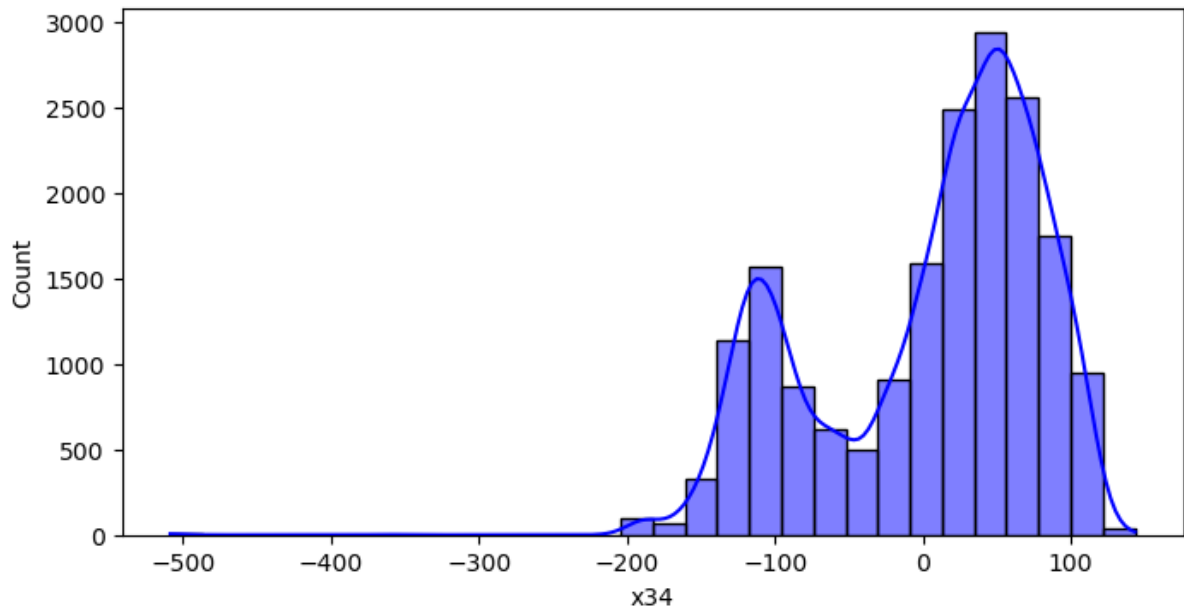




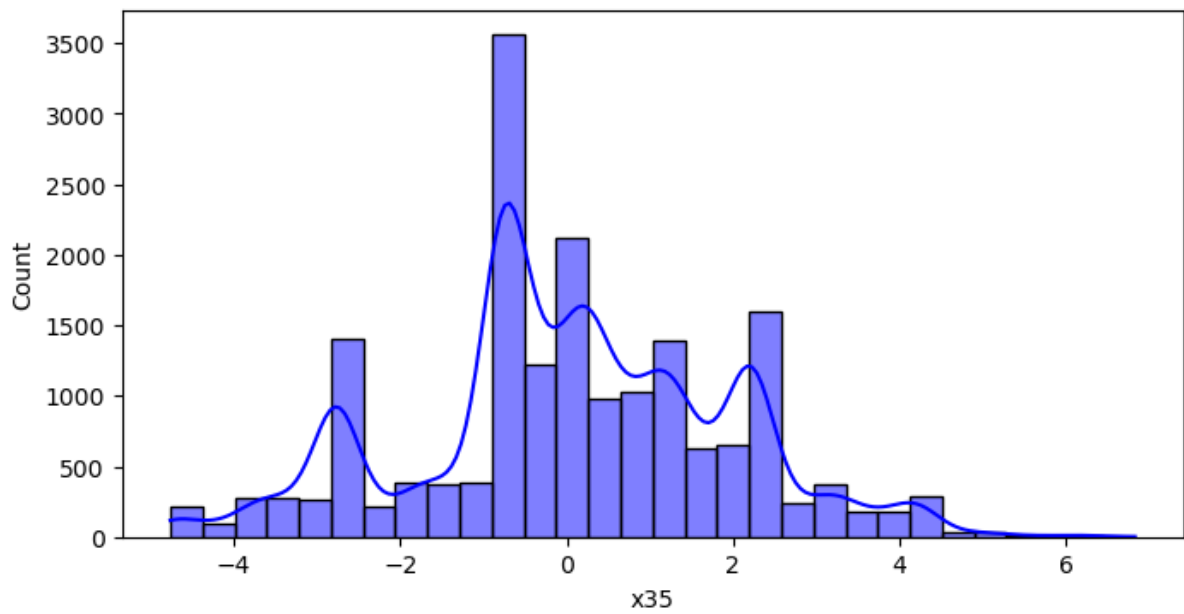




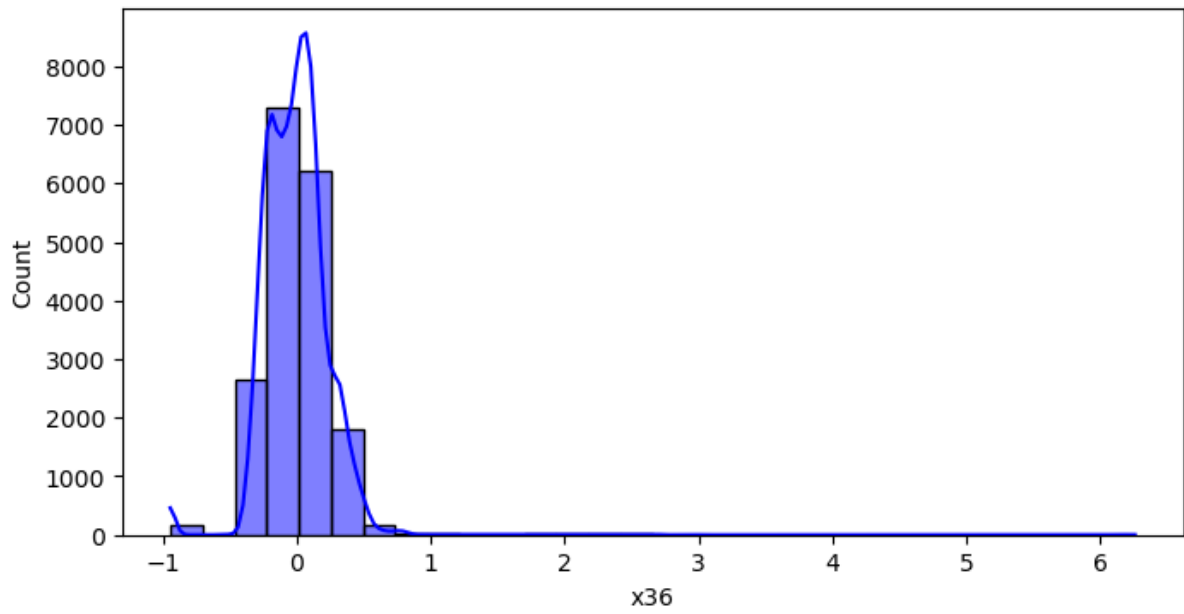
Distribution of x34



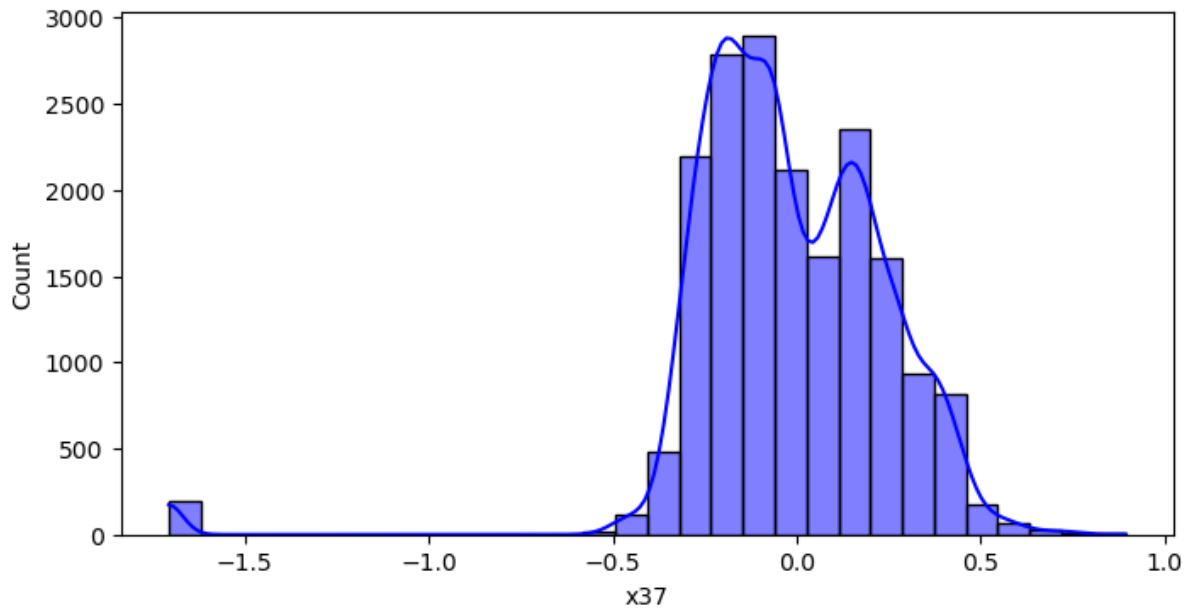
Distribution of x35



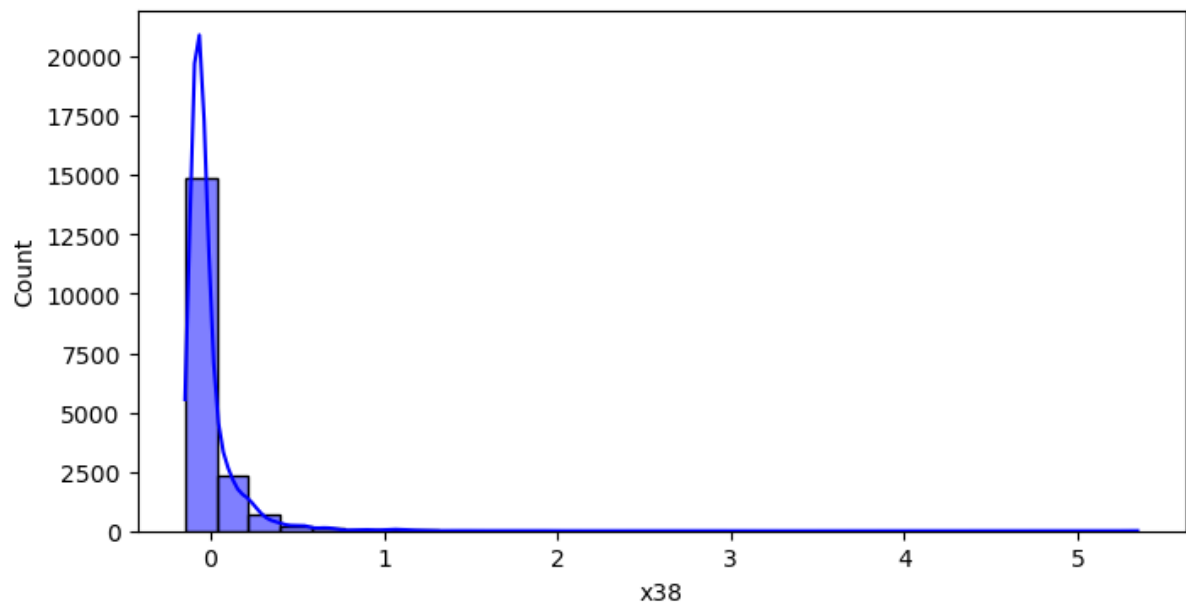
Distribution of x36



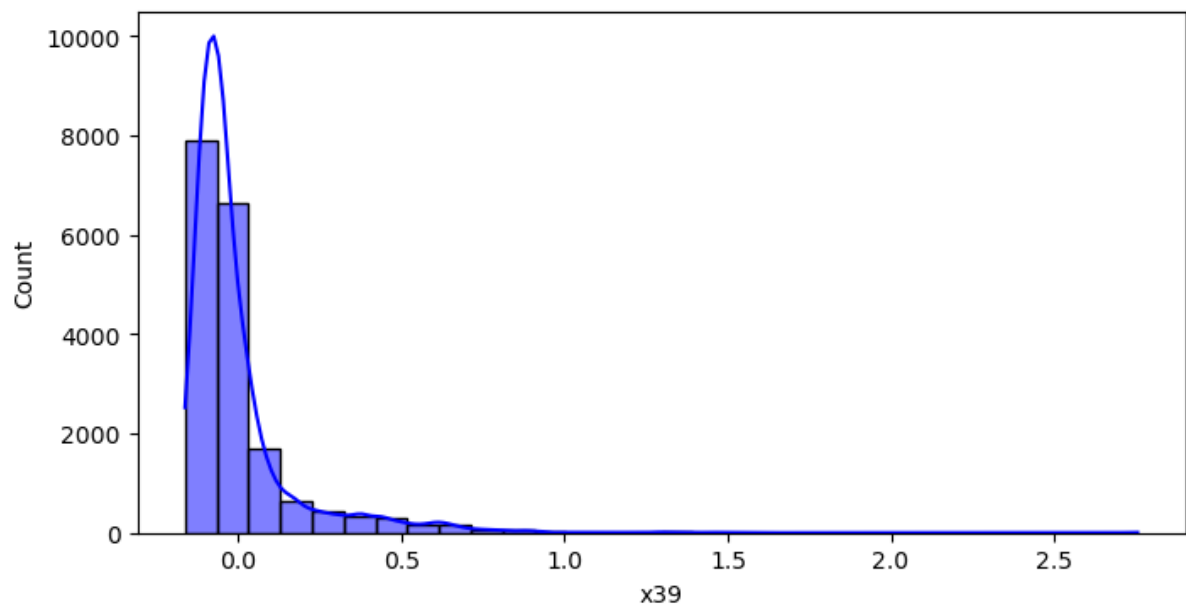
Distribution of x37

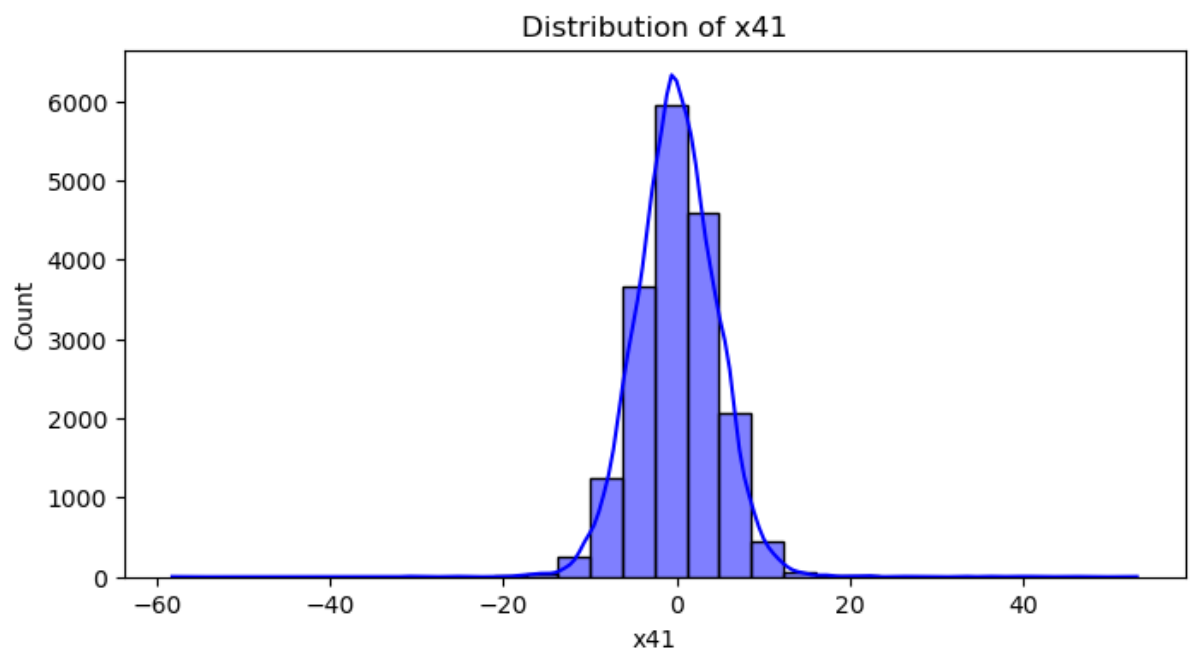
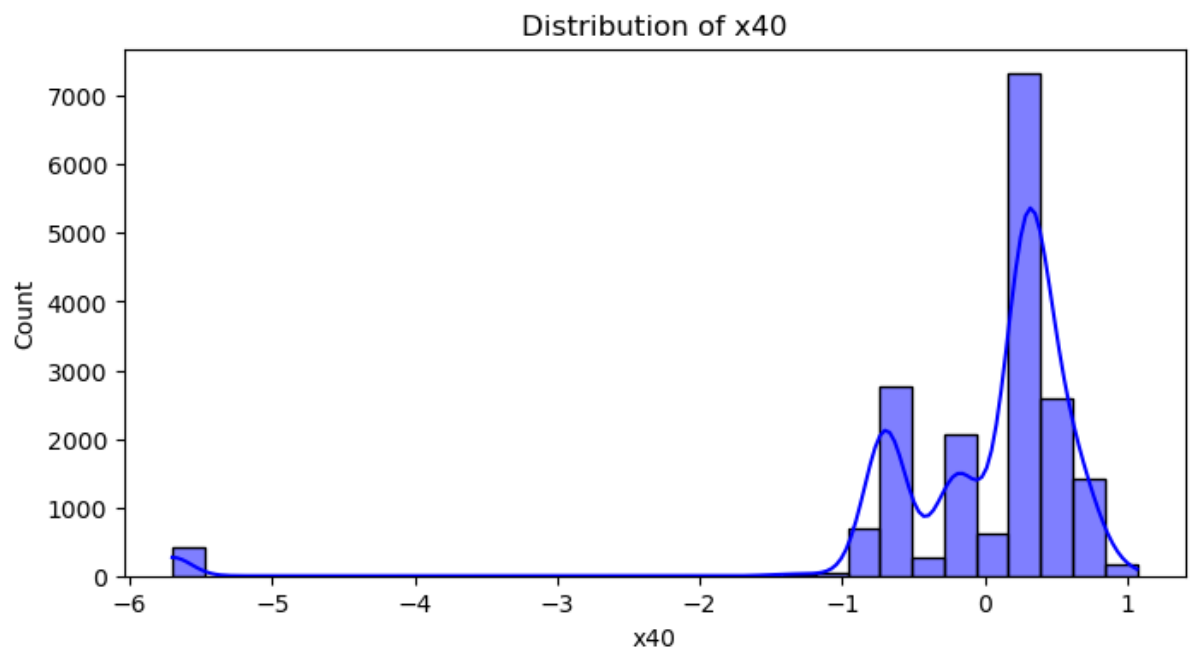


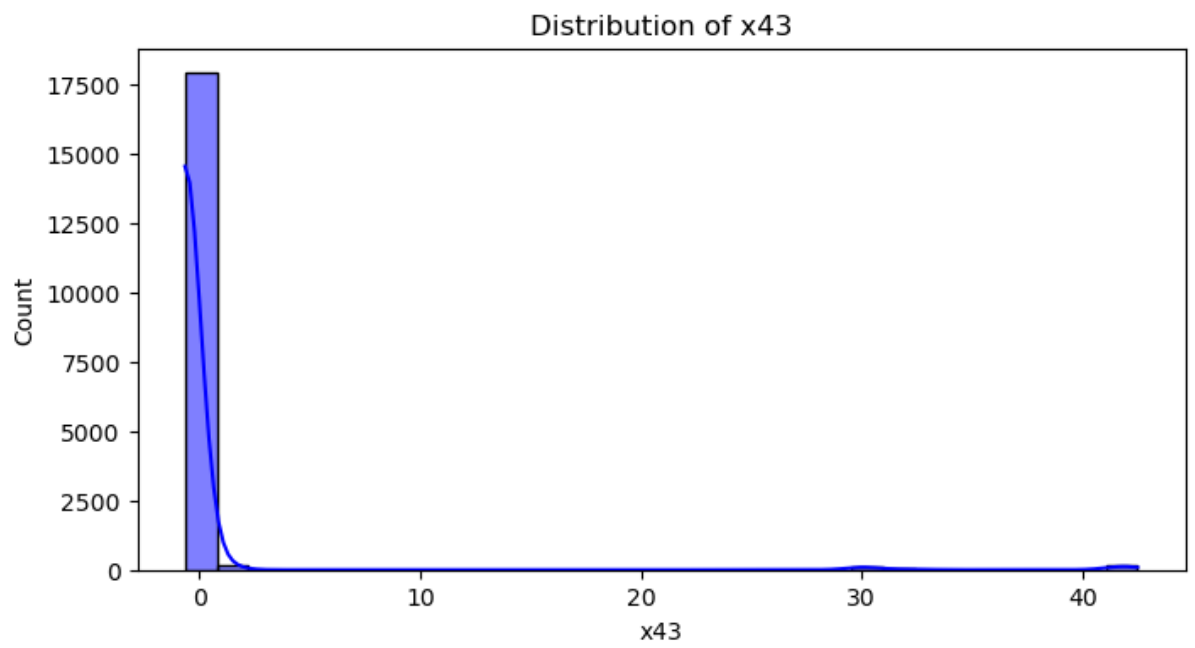
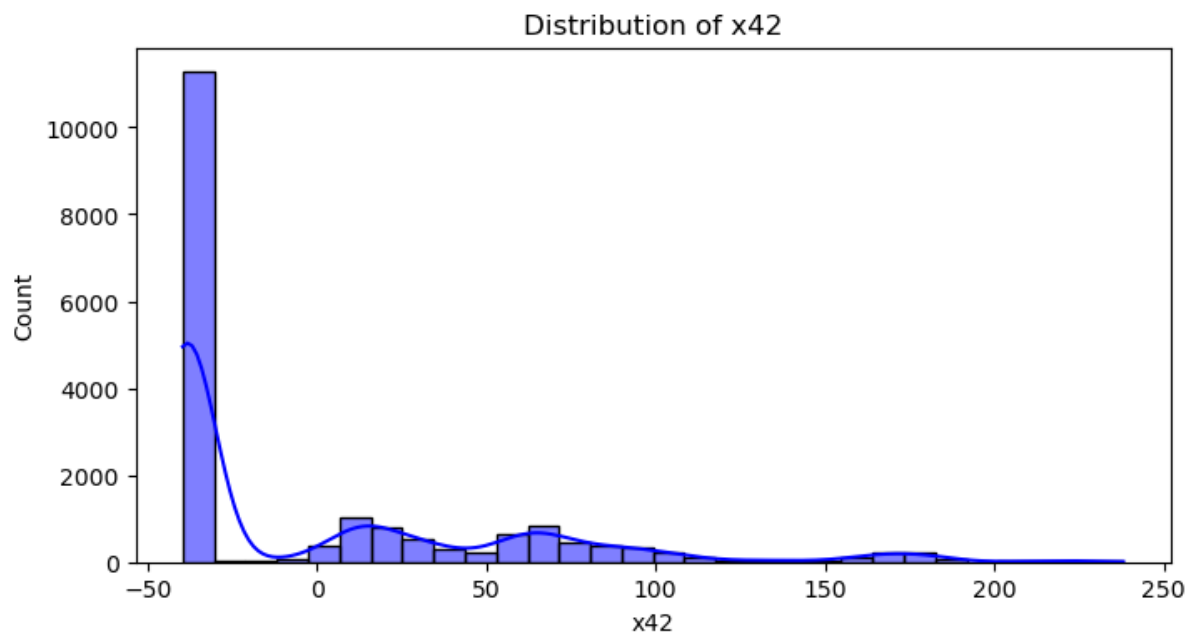
Distribution of x38

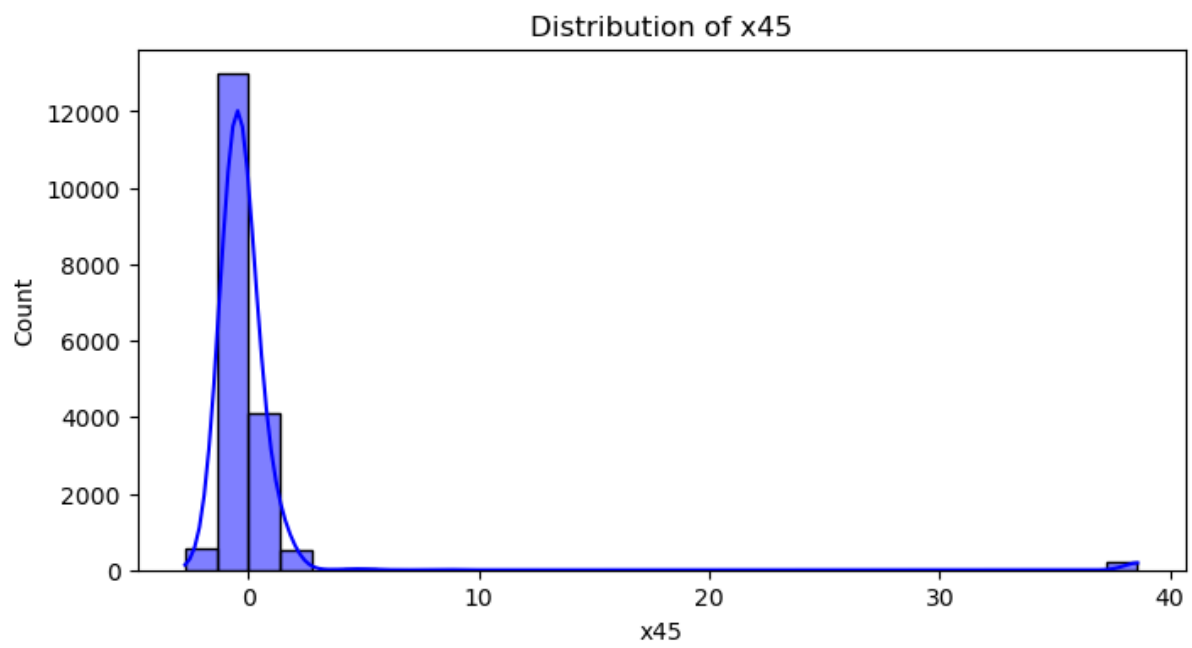
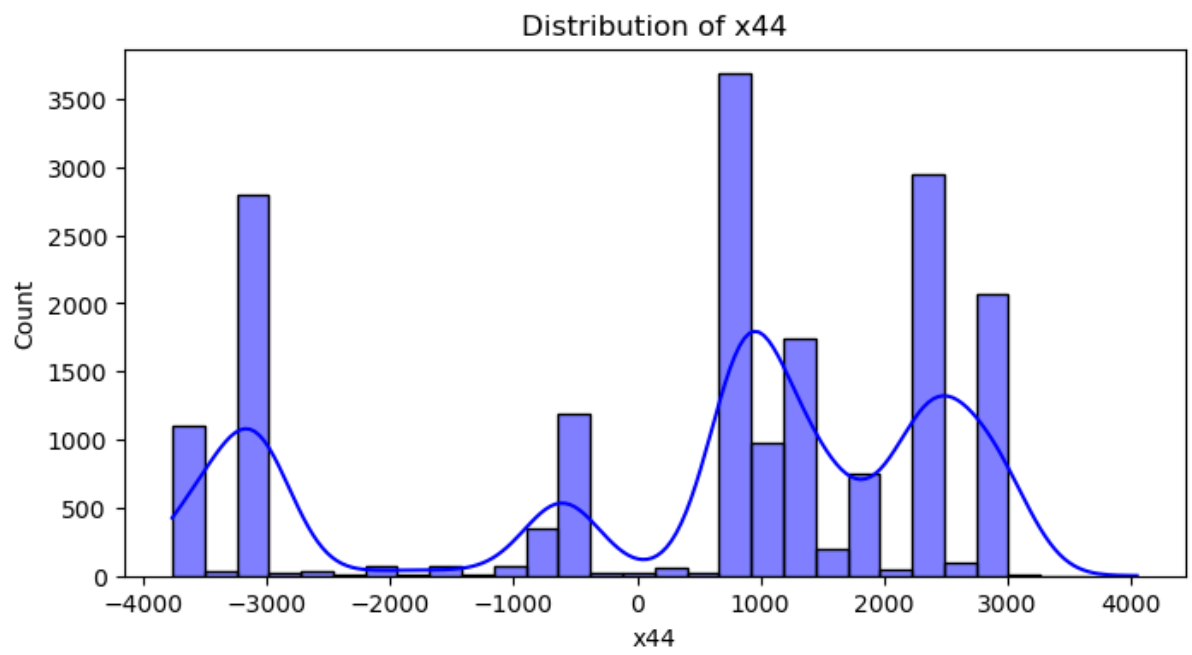


Distribution of x39

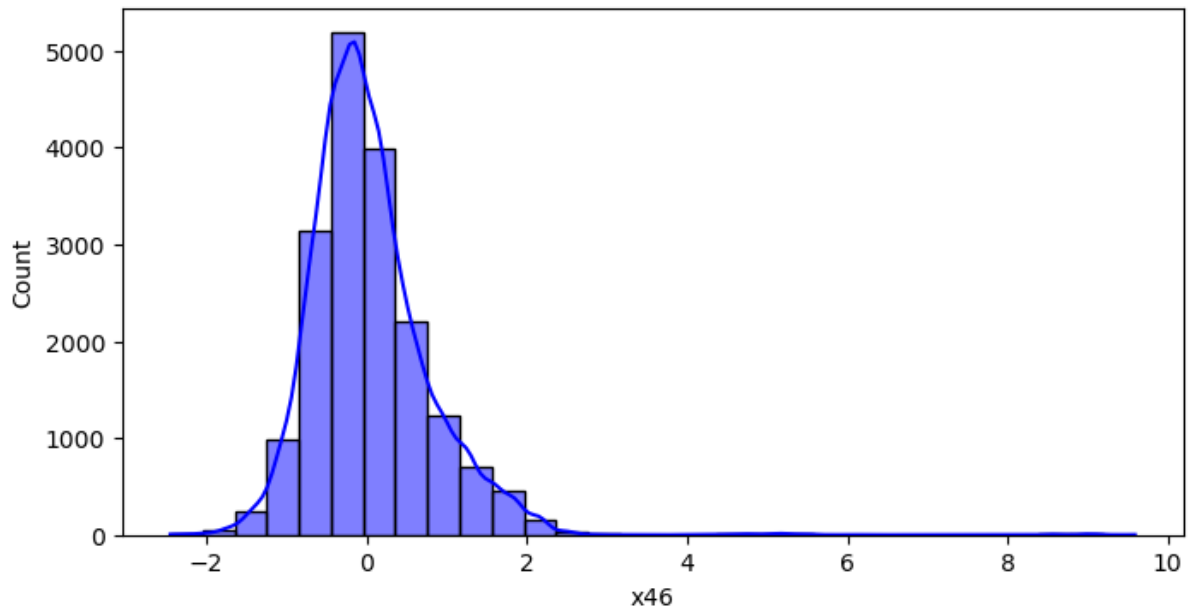




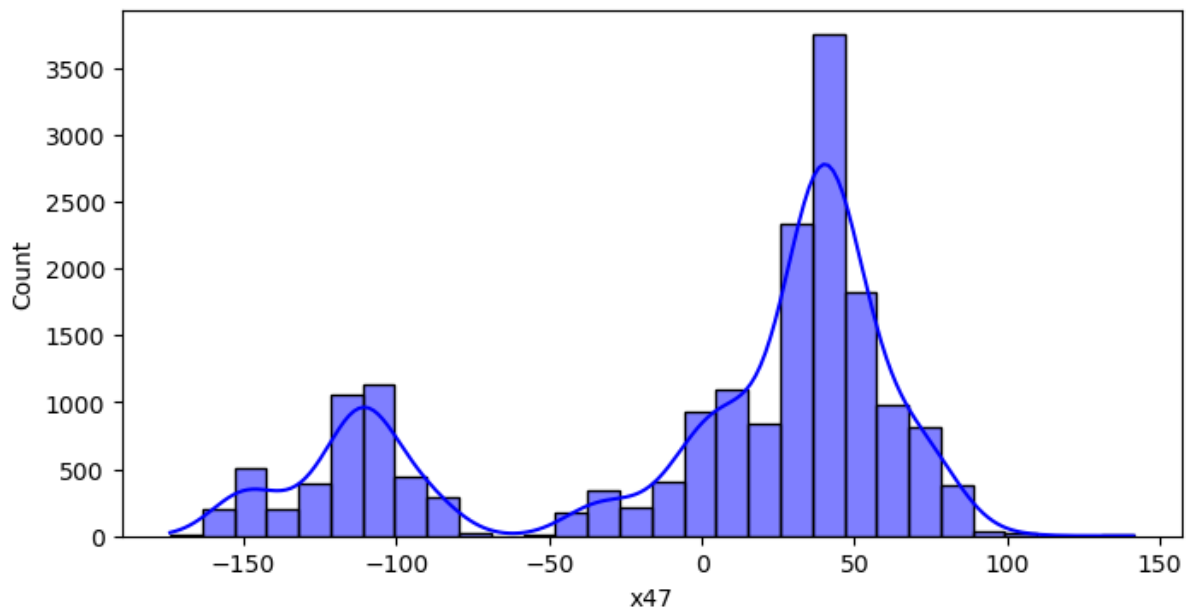


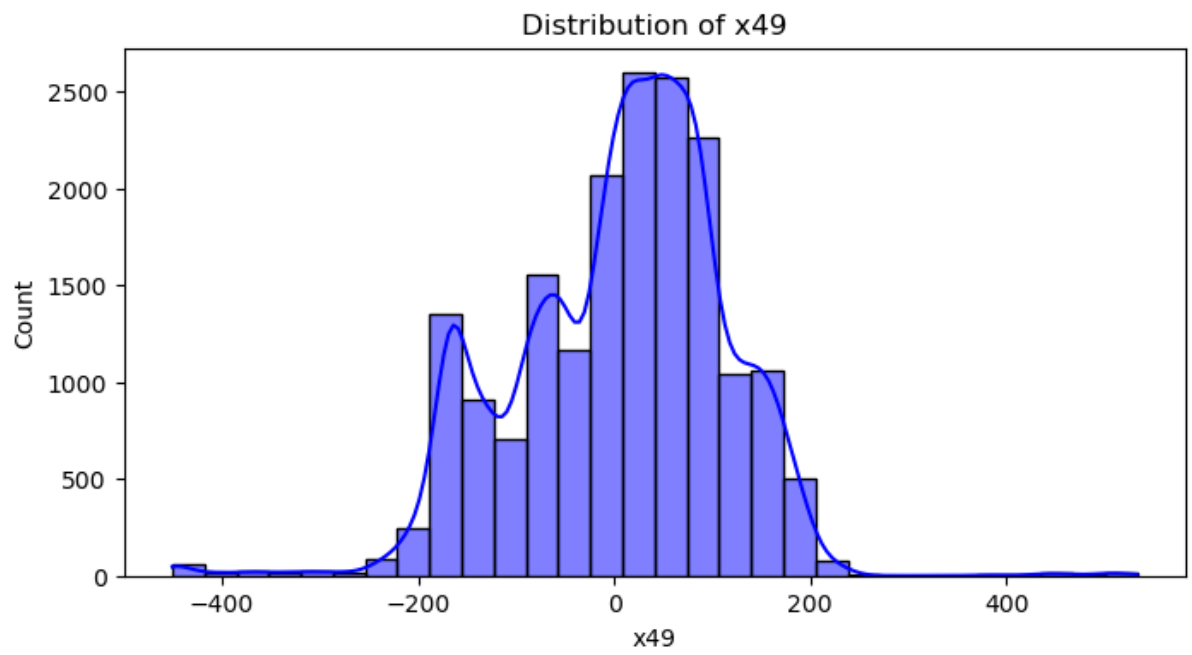
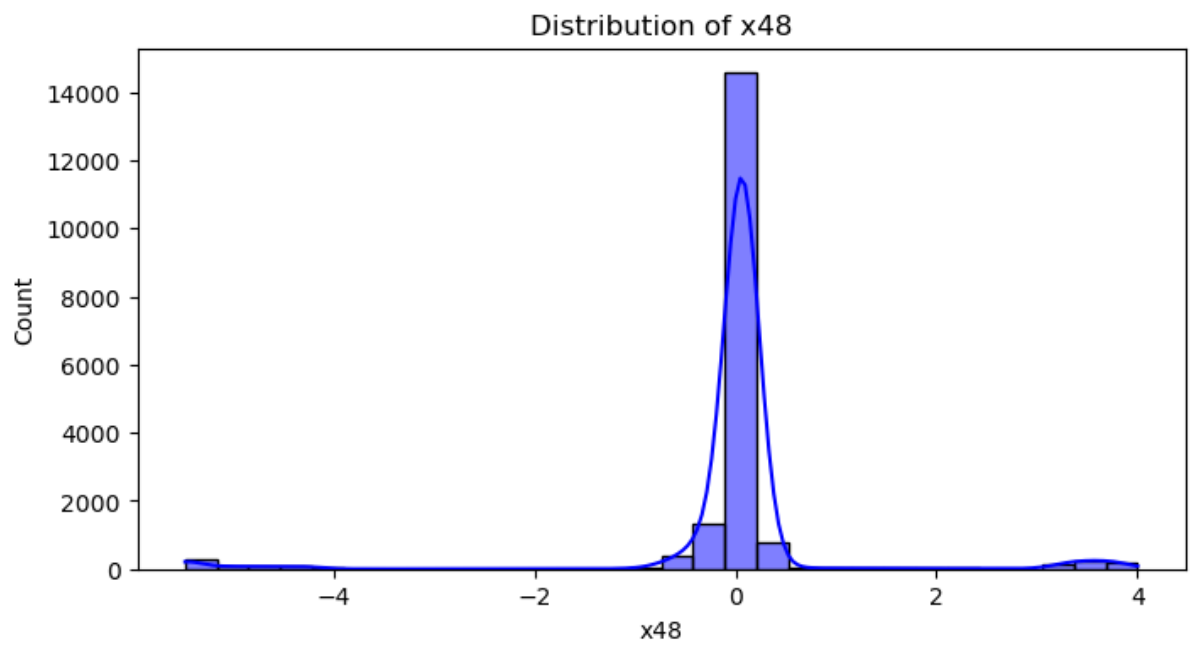


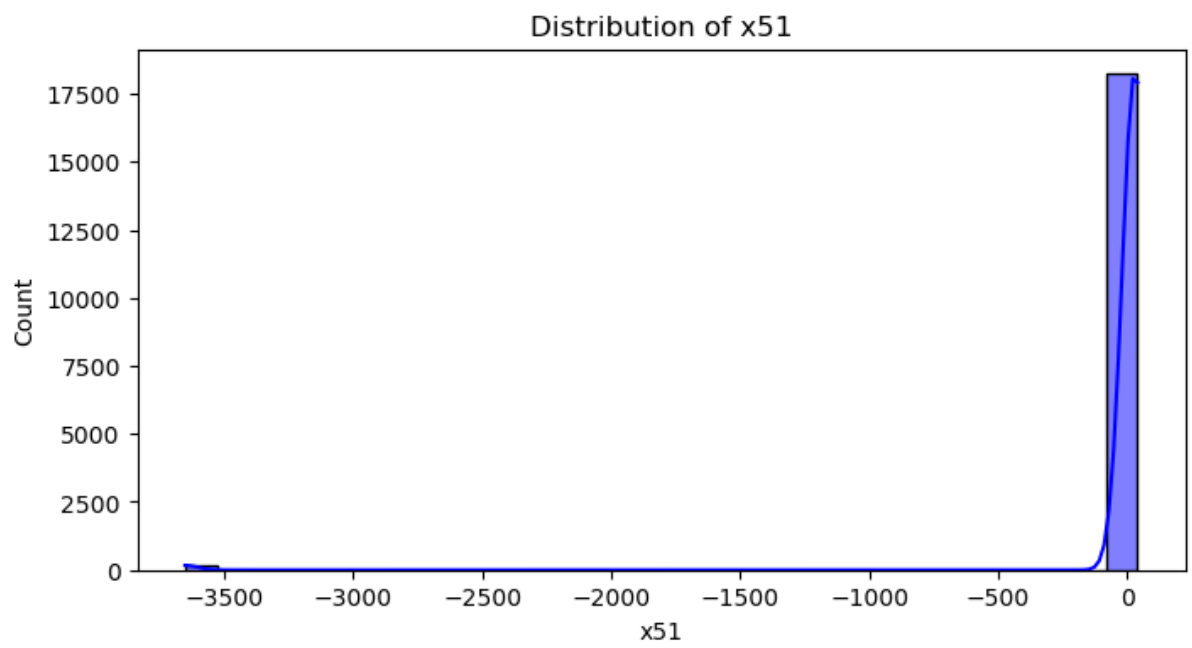
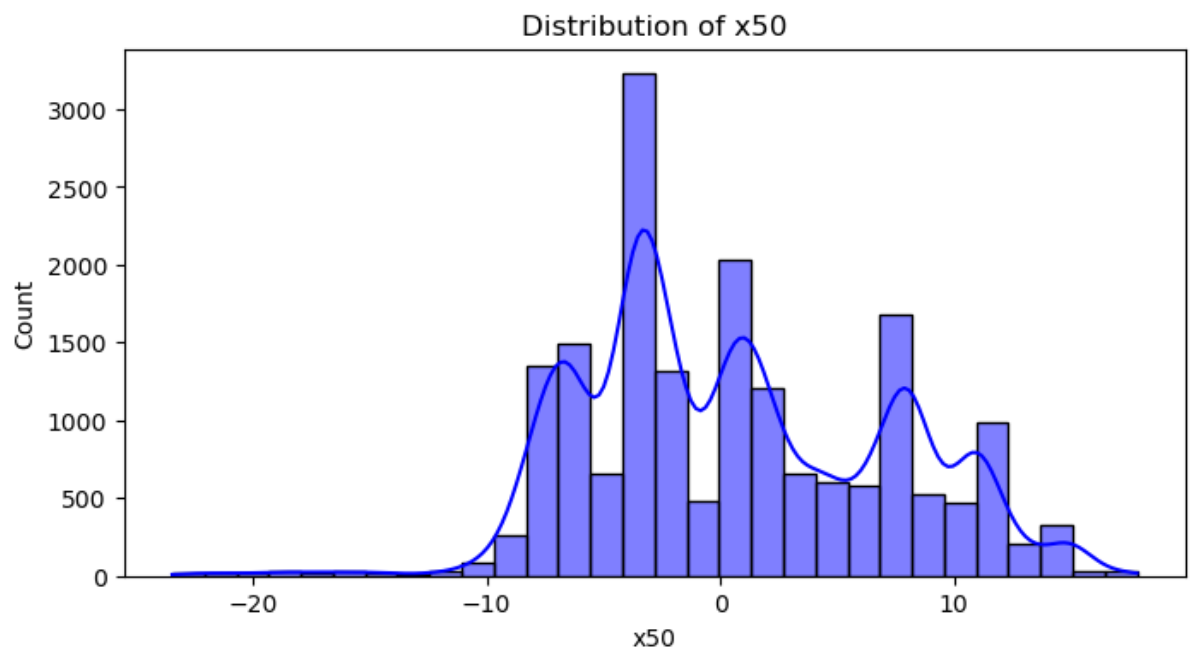
Distribution of x46

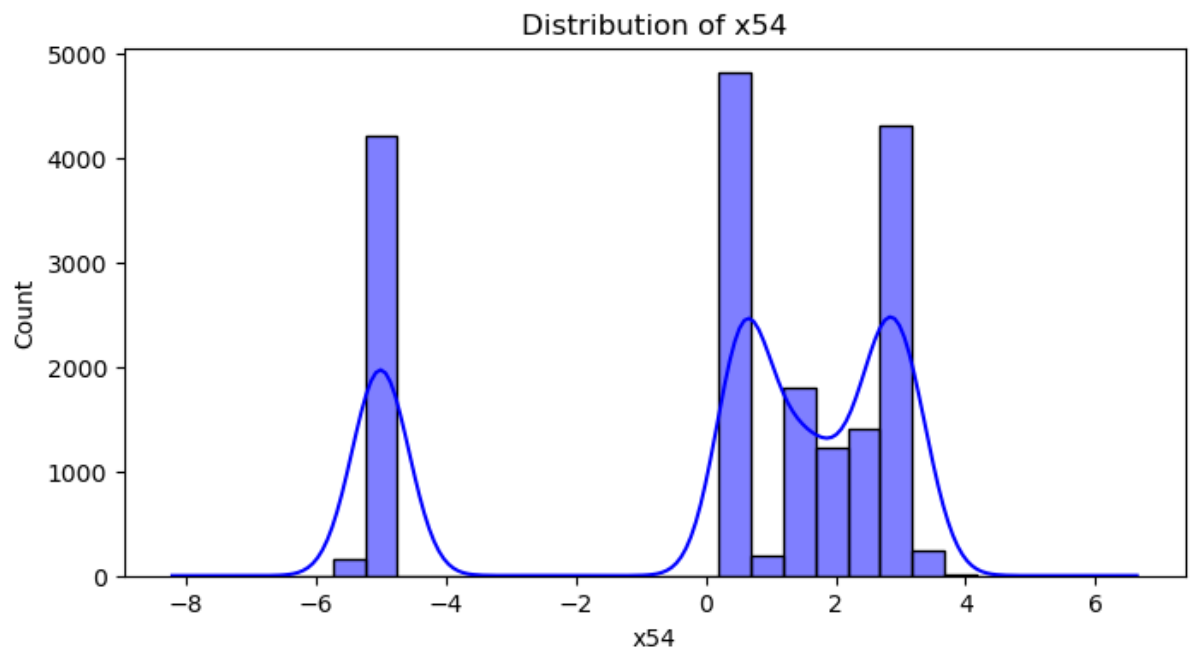
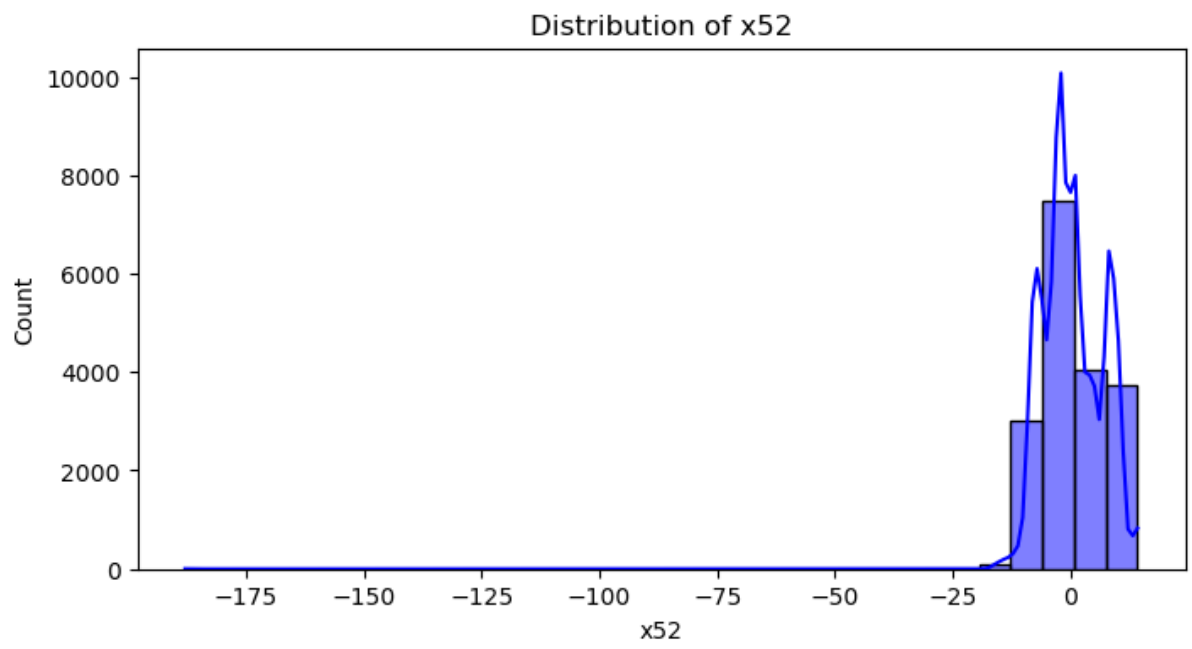


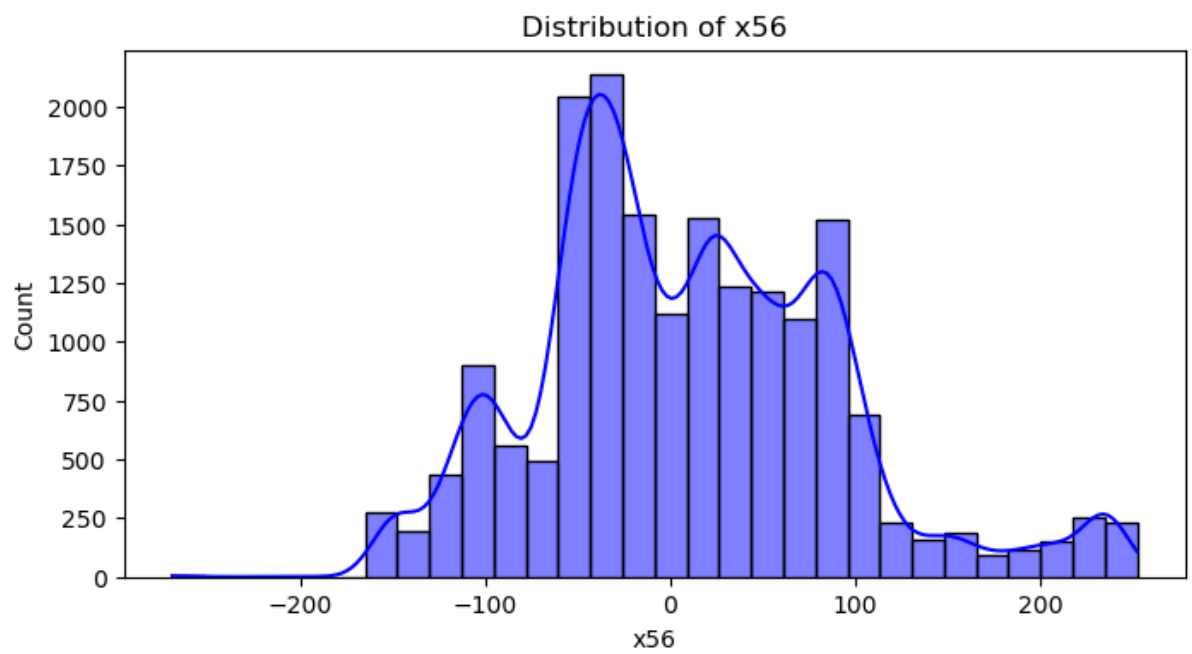
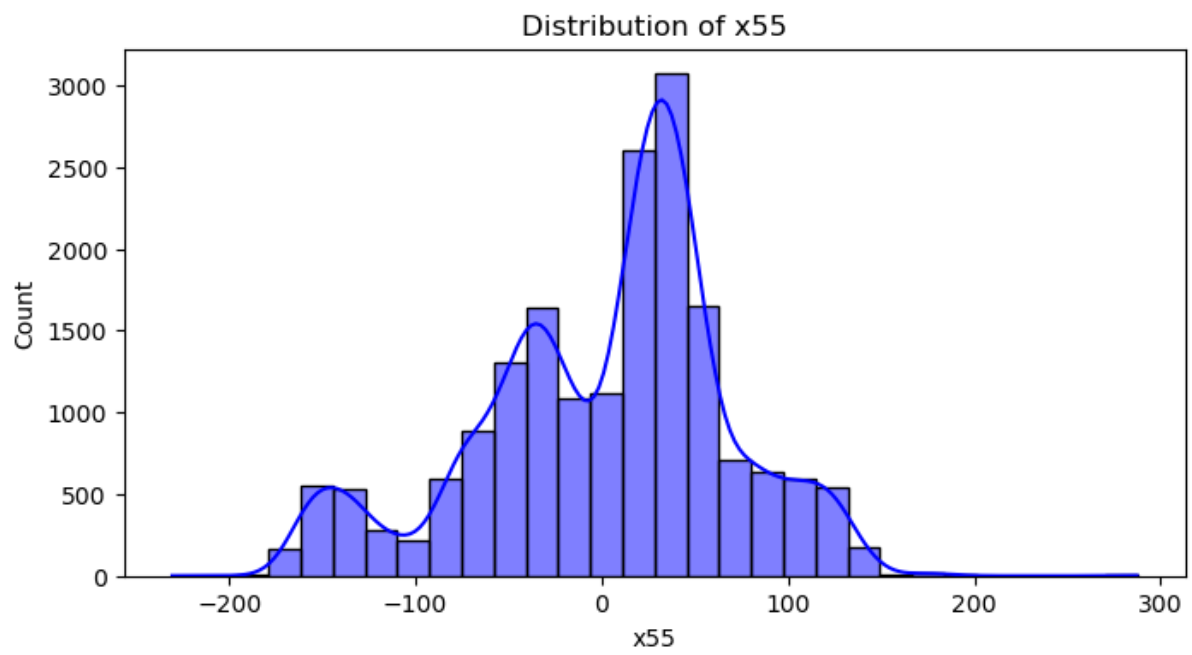
Distribution of x47



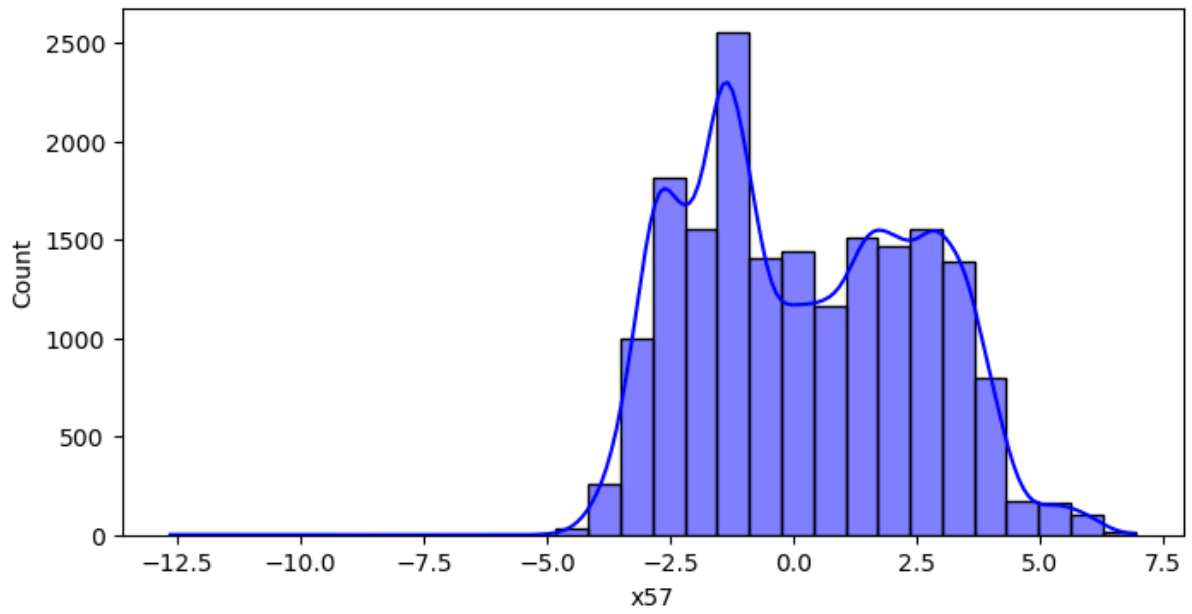




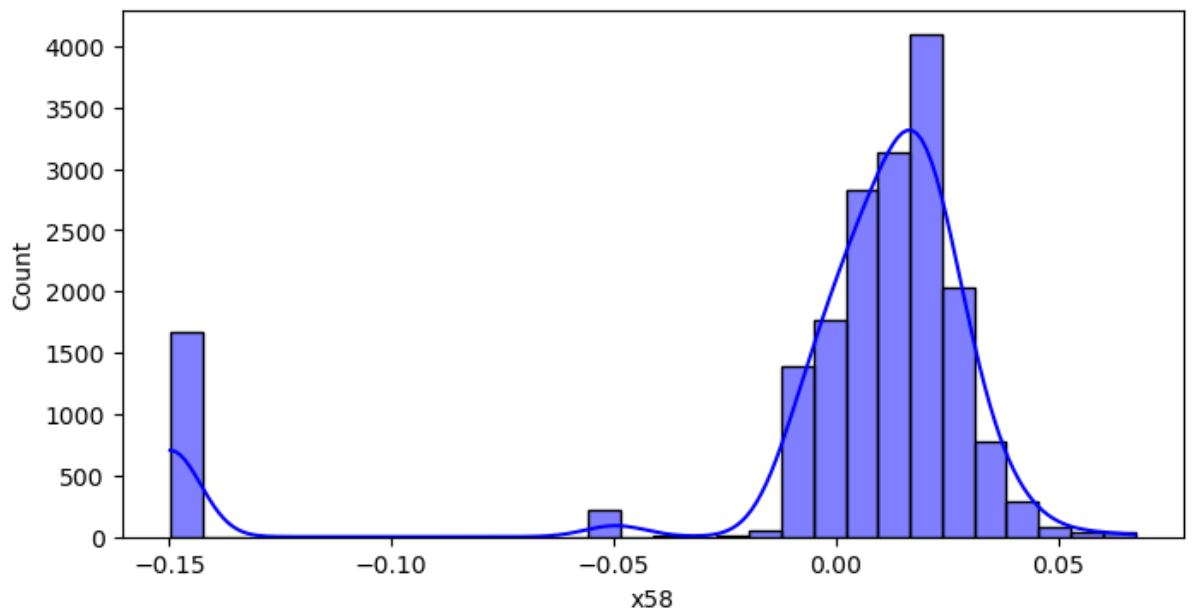


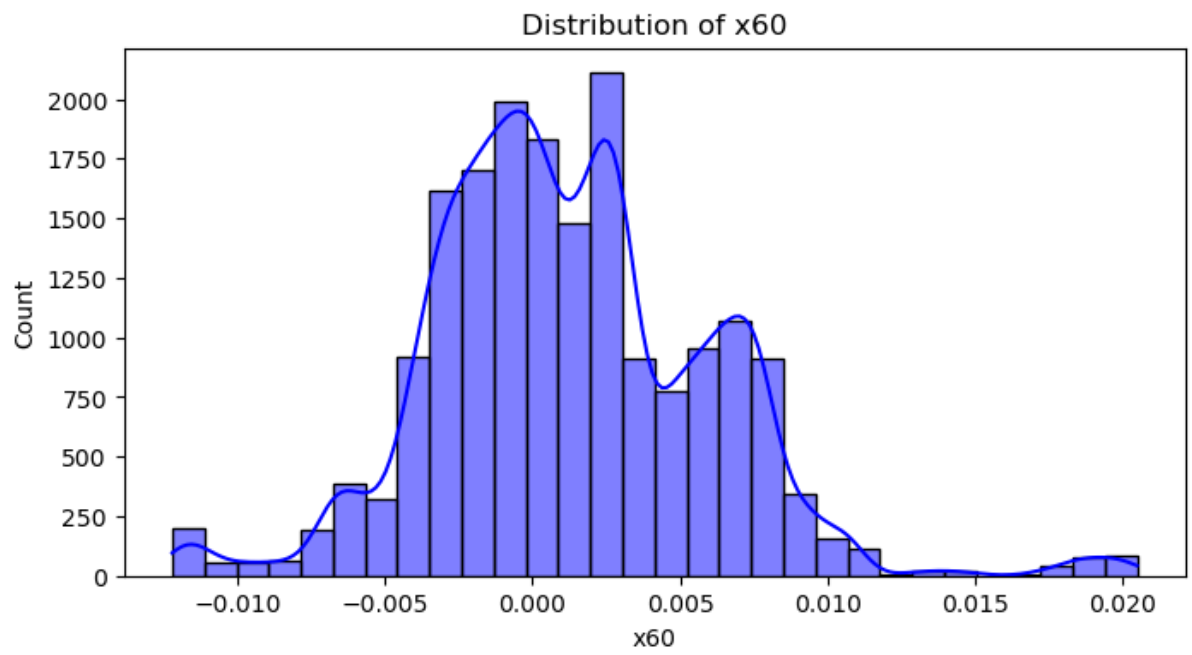
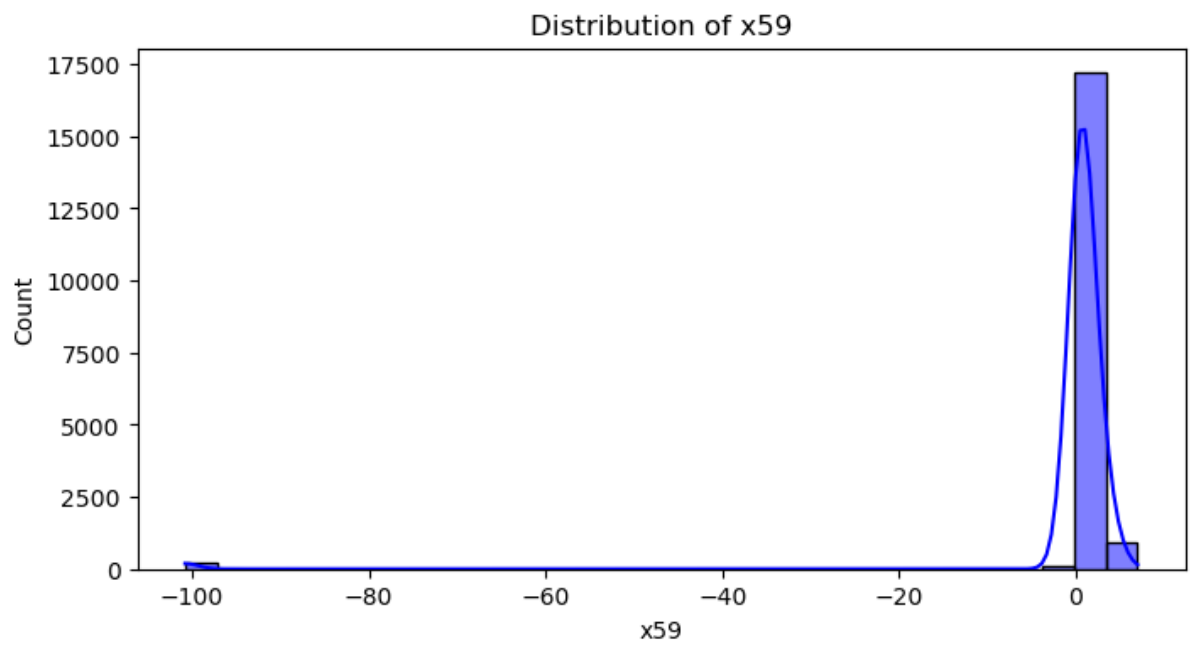


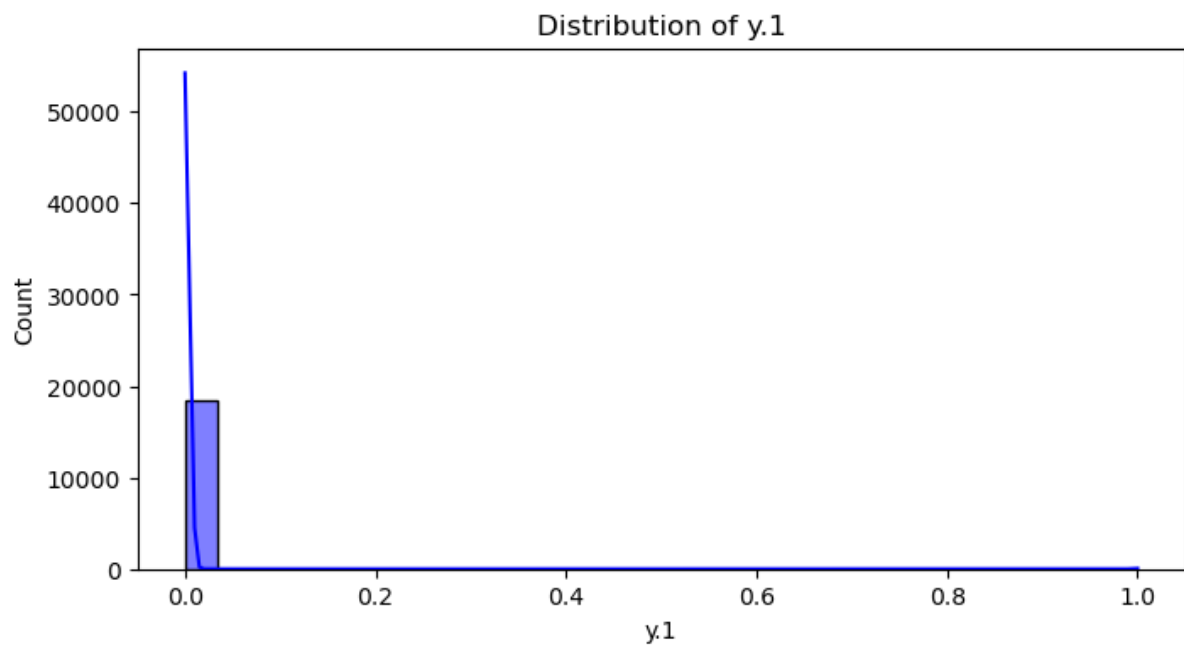
Distribution of x57



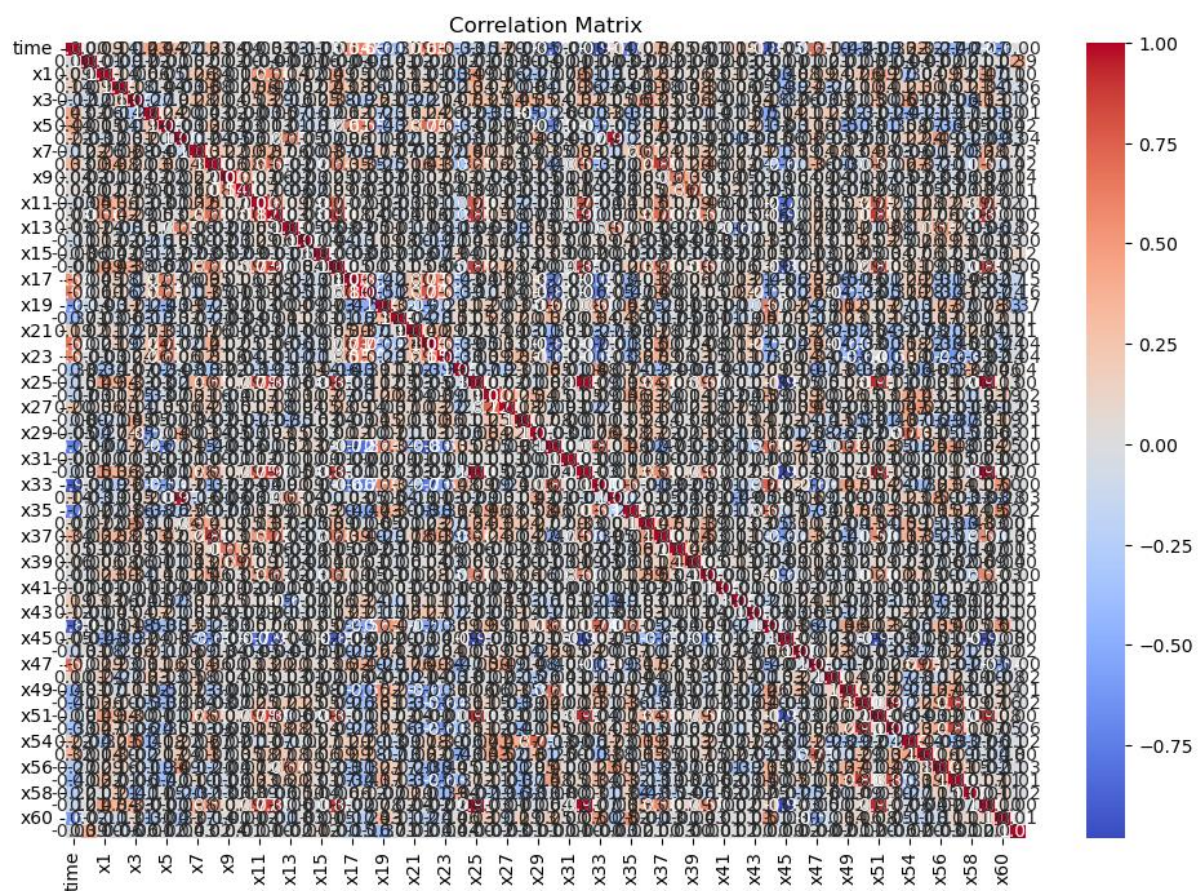
Distribution of x58





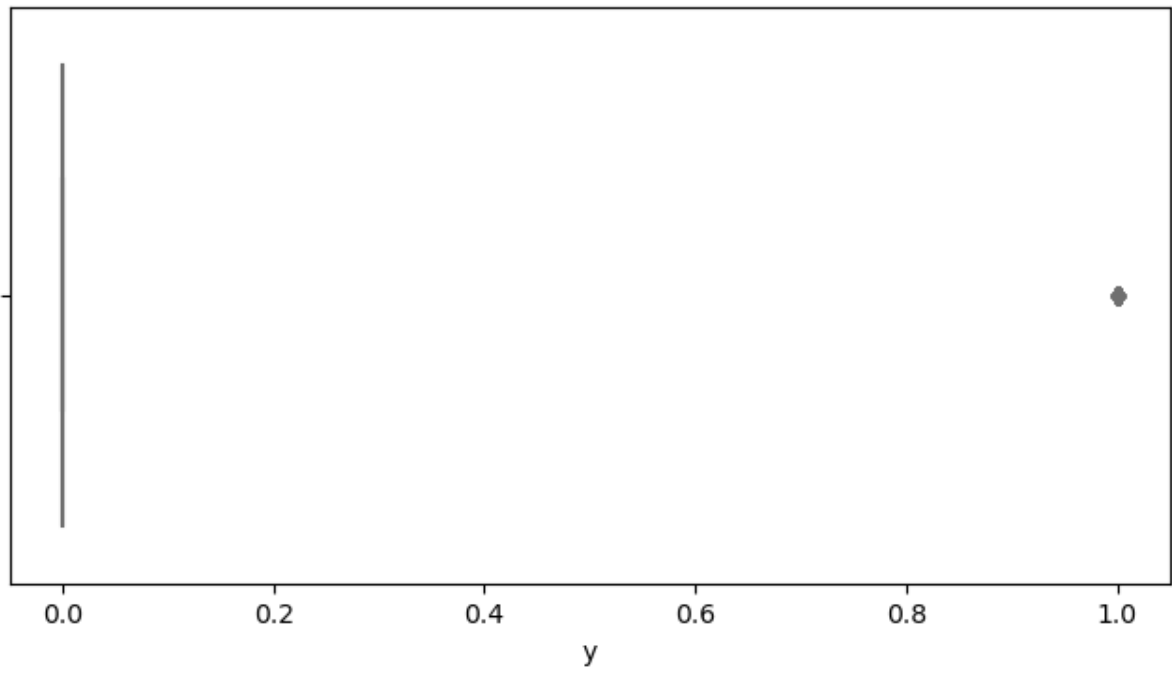


Correlation Heatmap:

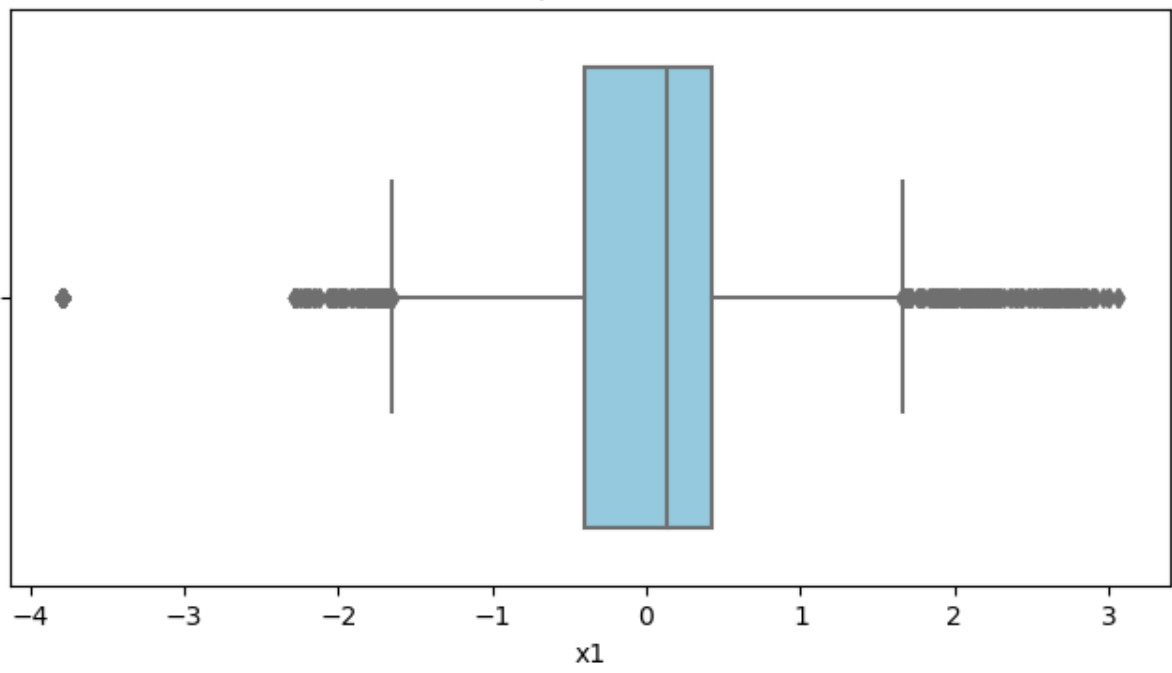


Generating boxplots for numerical features:

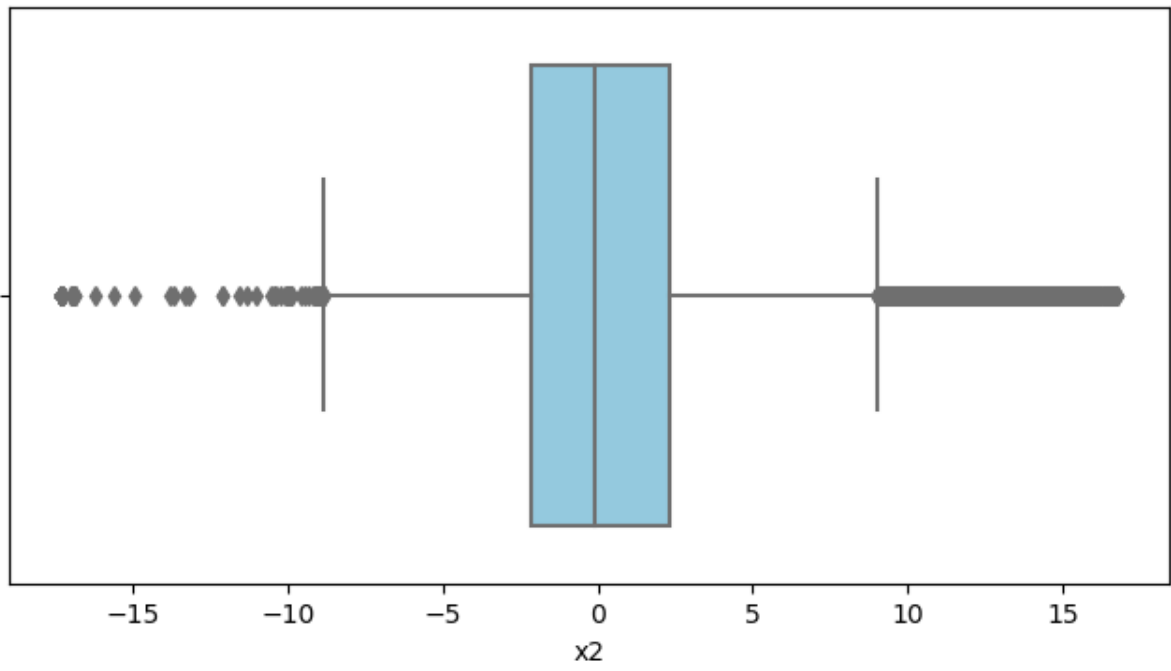
Boxplot for y



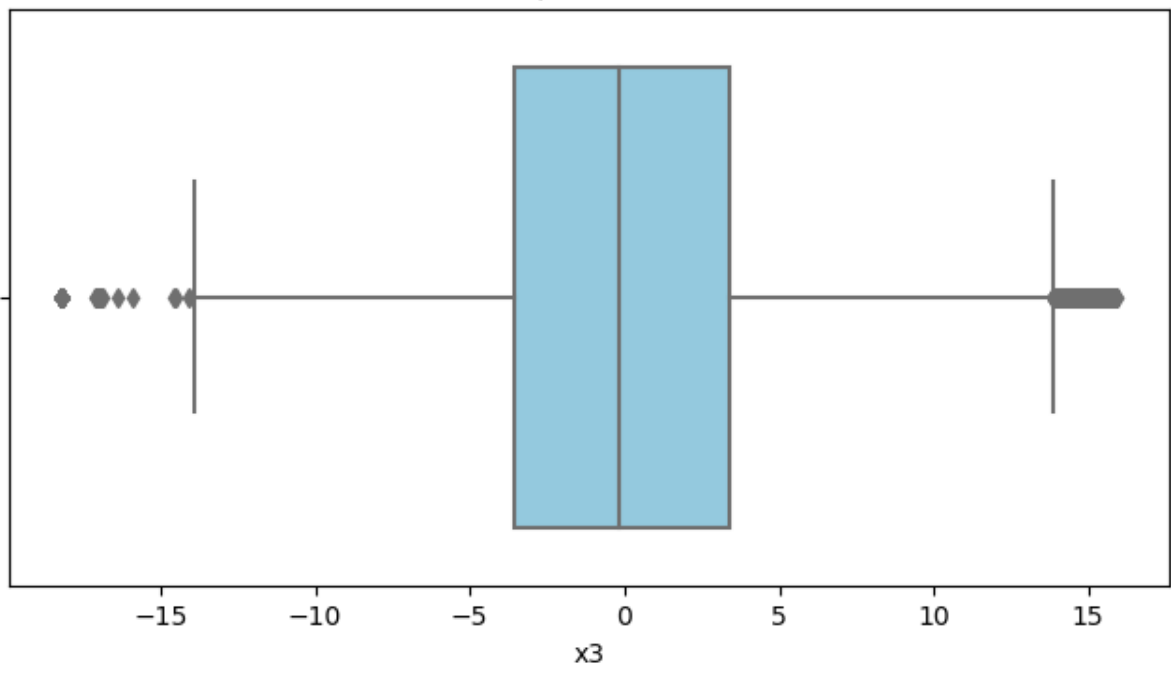
Boxplot for x1



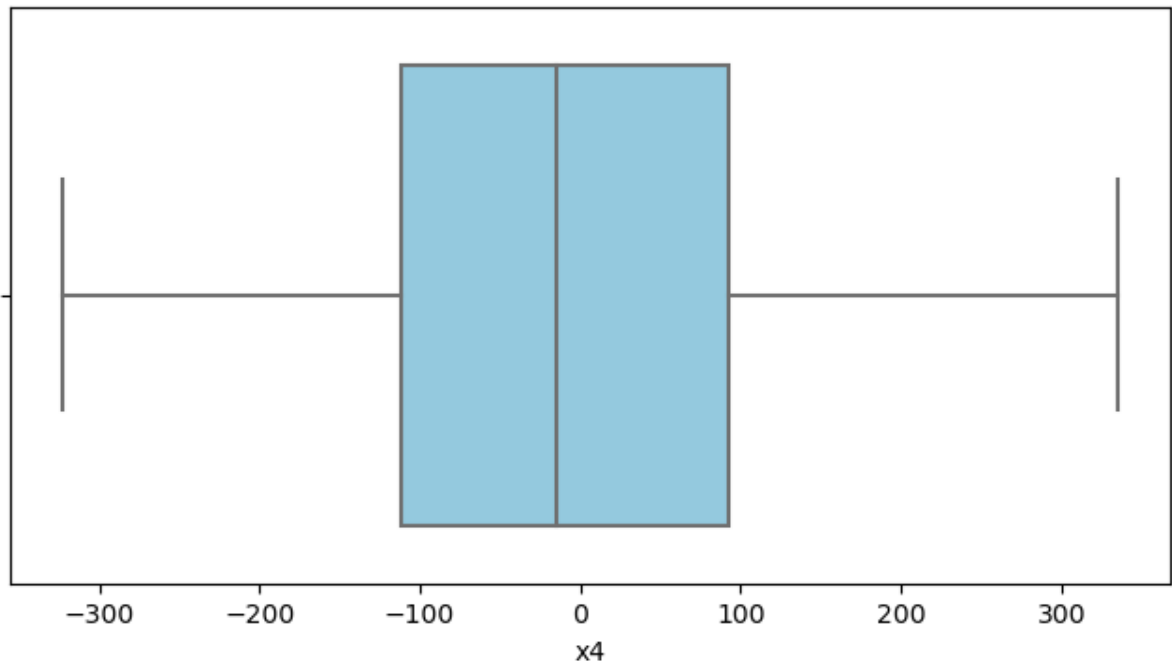
Boxplot for x2



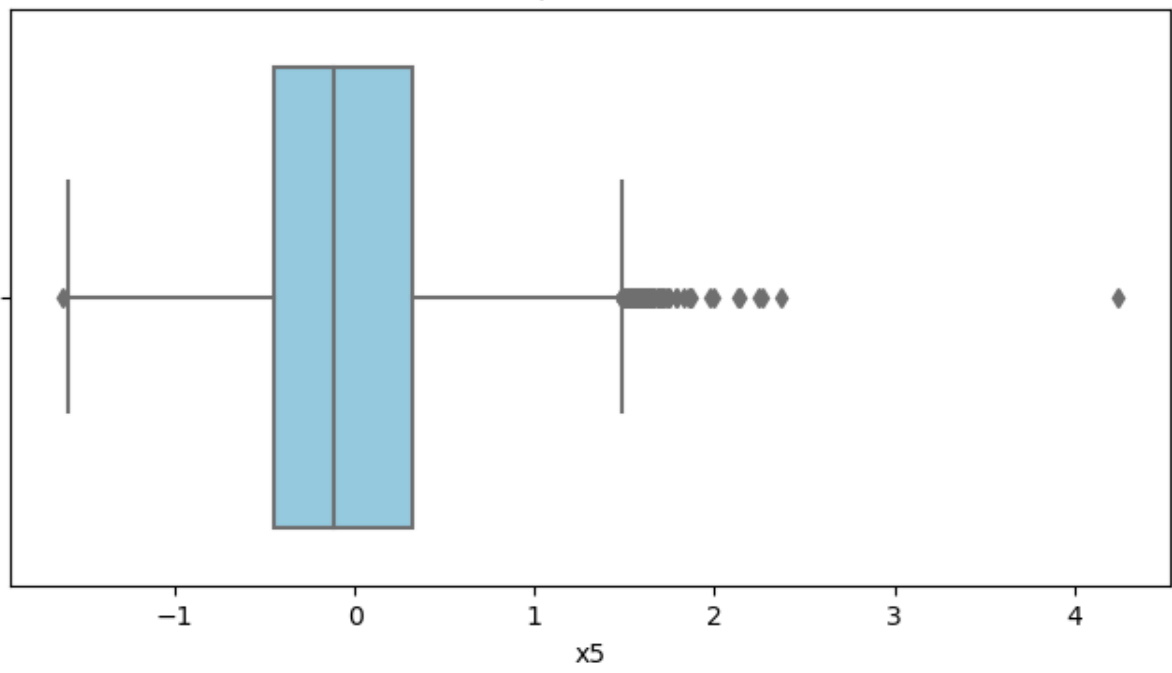
Boxplot for x3



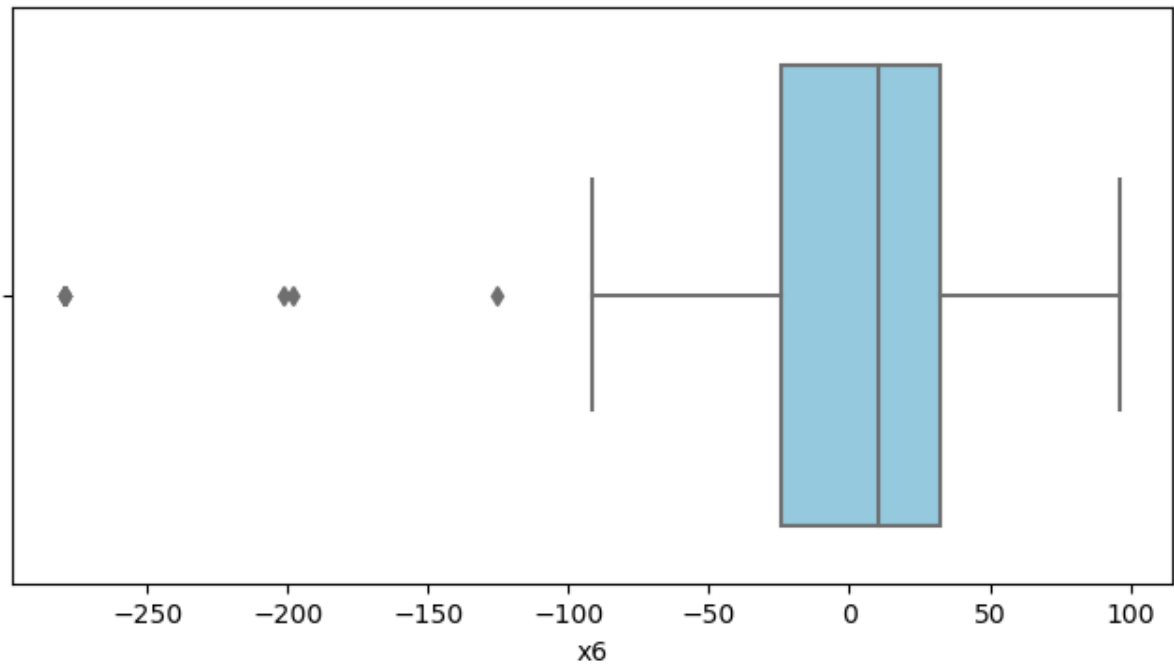
Boxplot for x4



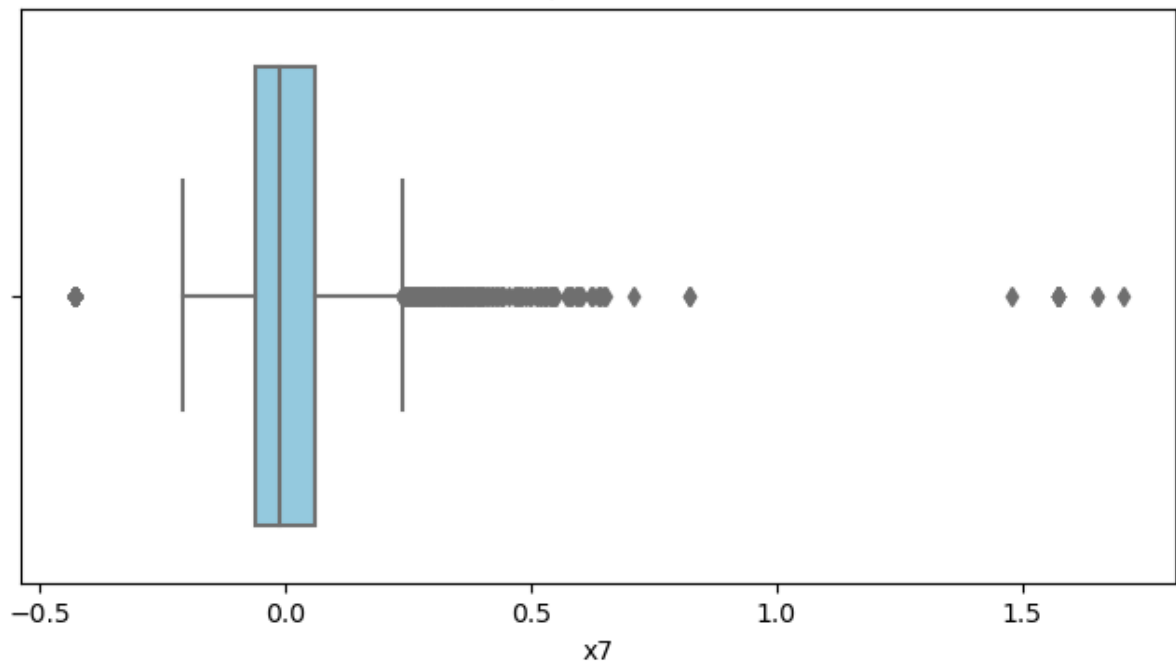
Boxplot for x5



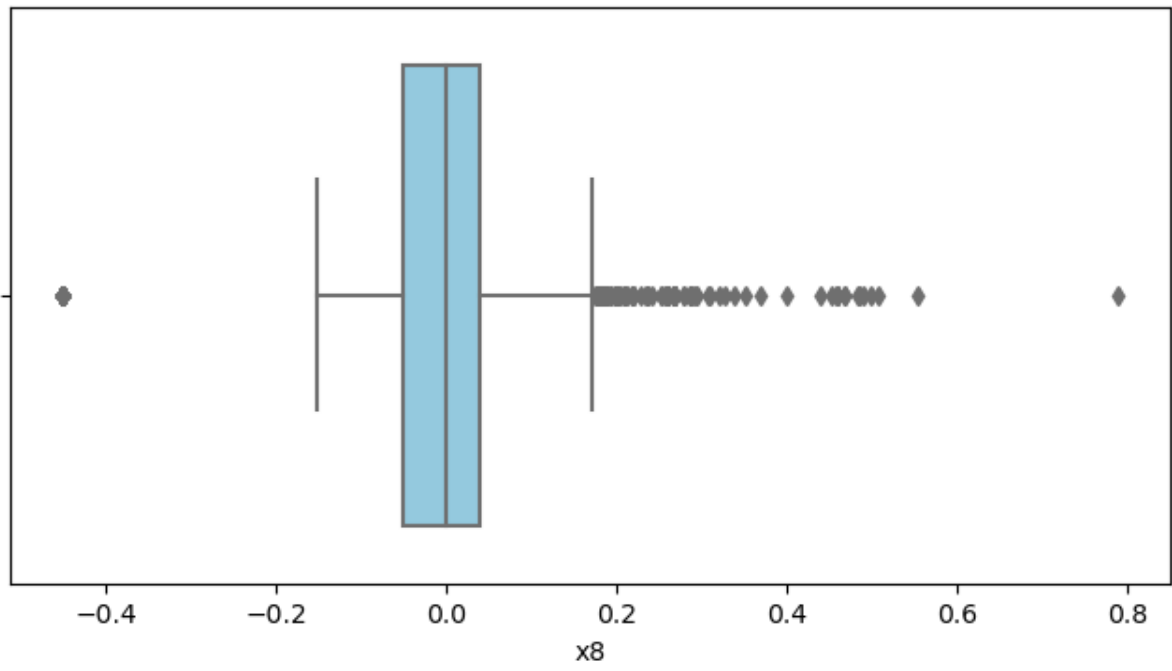
Boxplot for x6



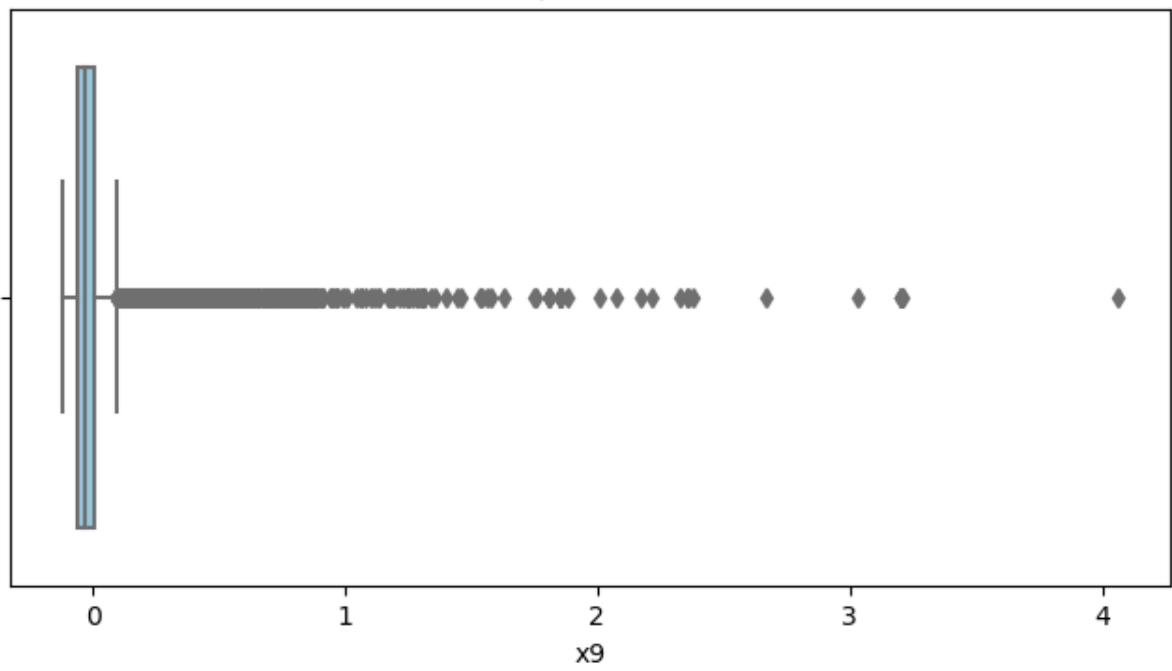
Boxplot for x7



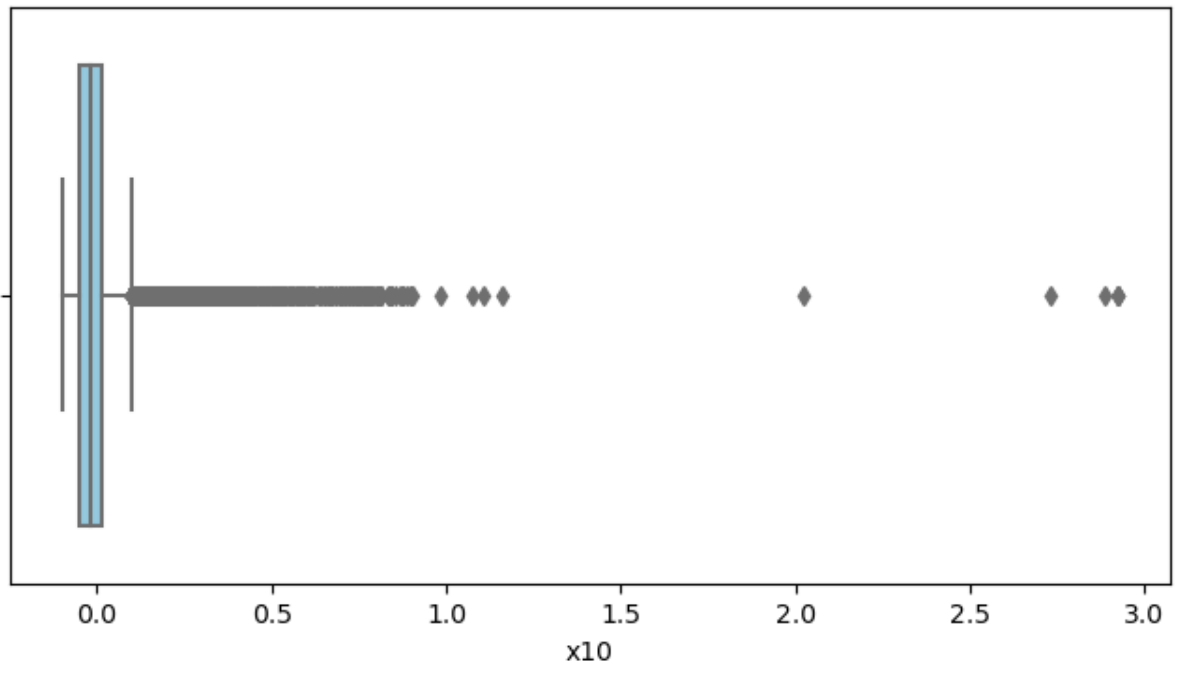
Boxplot for x8



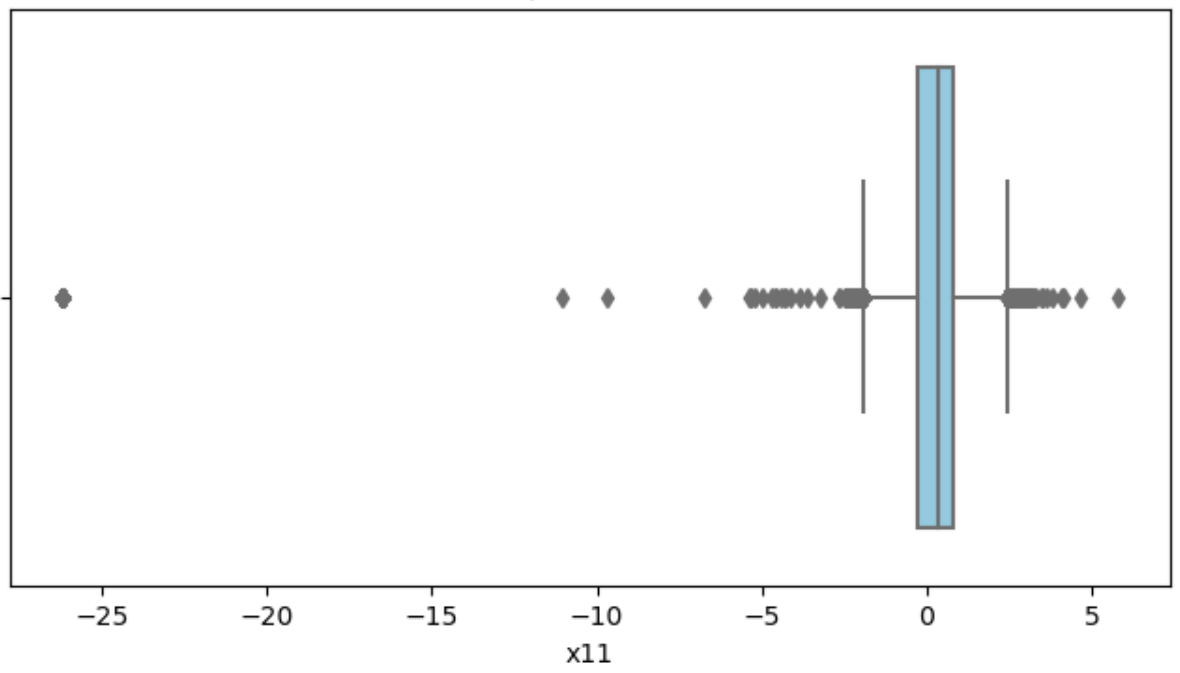
Boxplot for x9



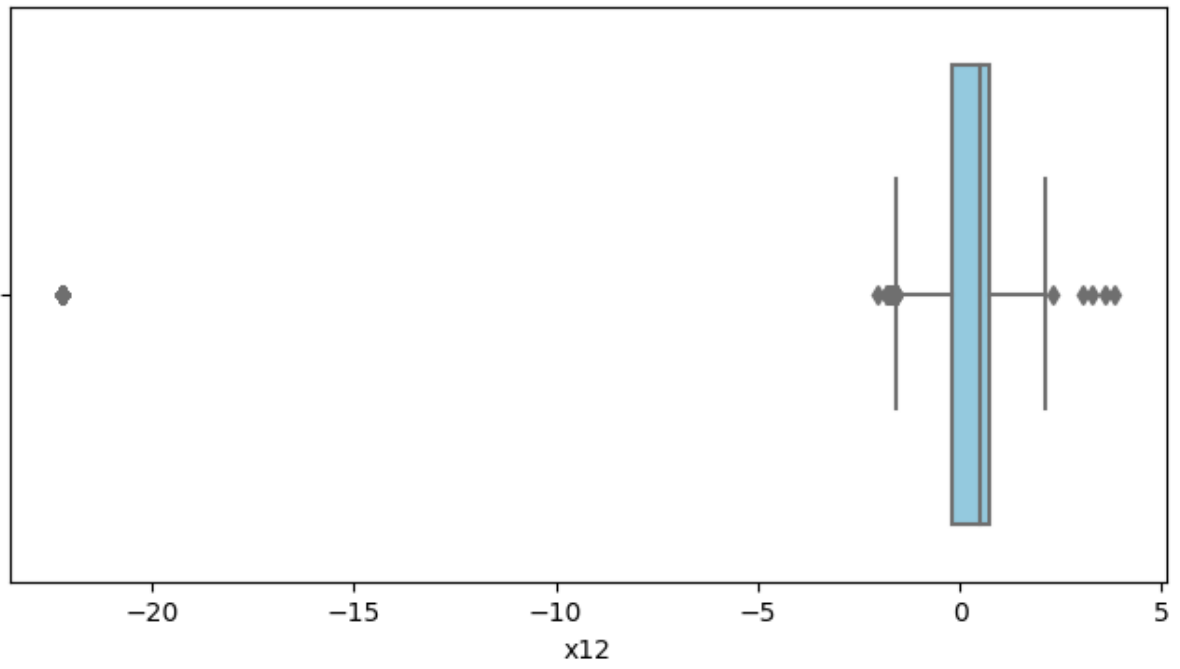
Boxplot for x10



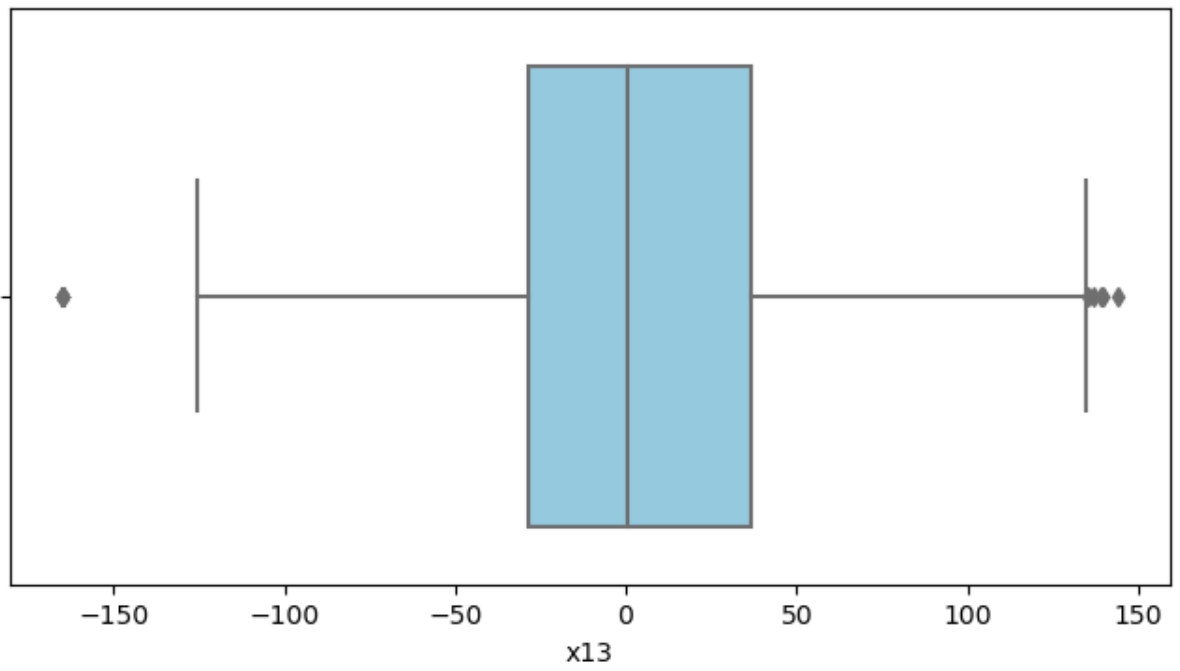
Boxplot for x11



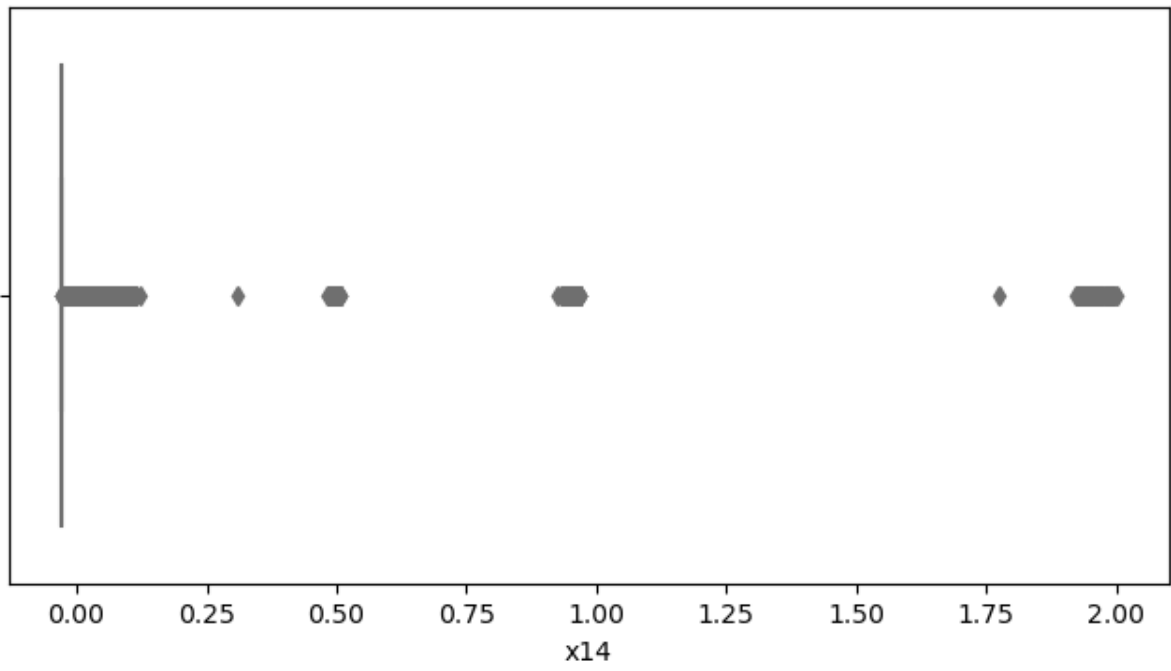
Boxplot for x12



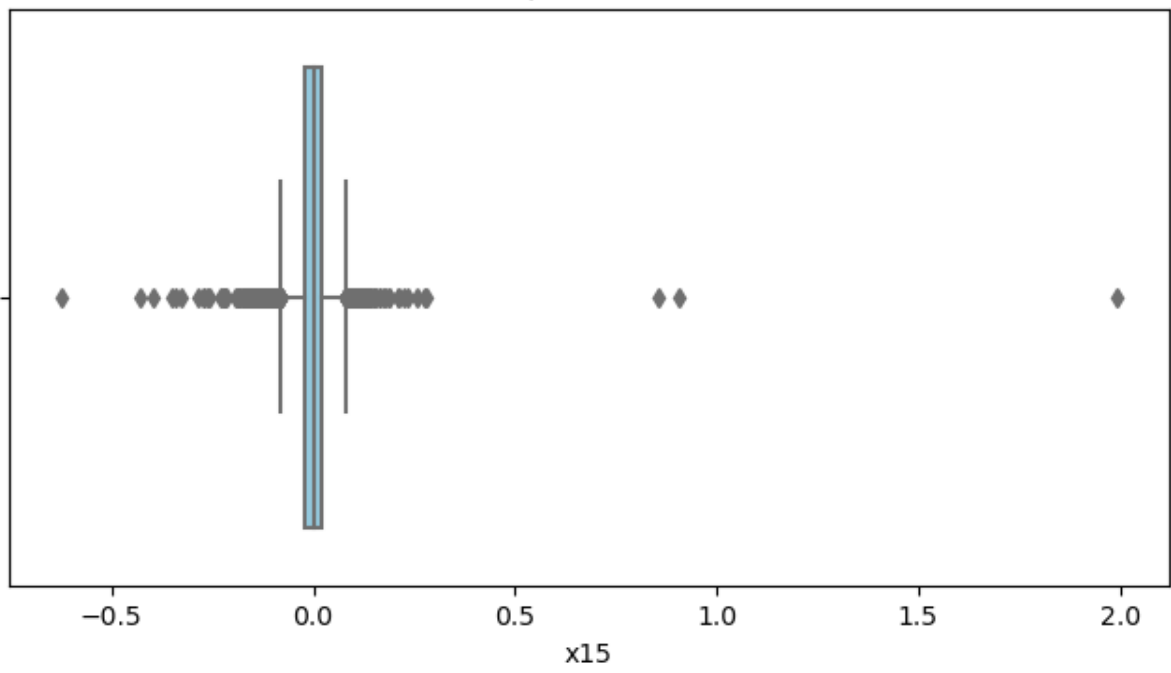
Boxplot for x13



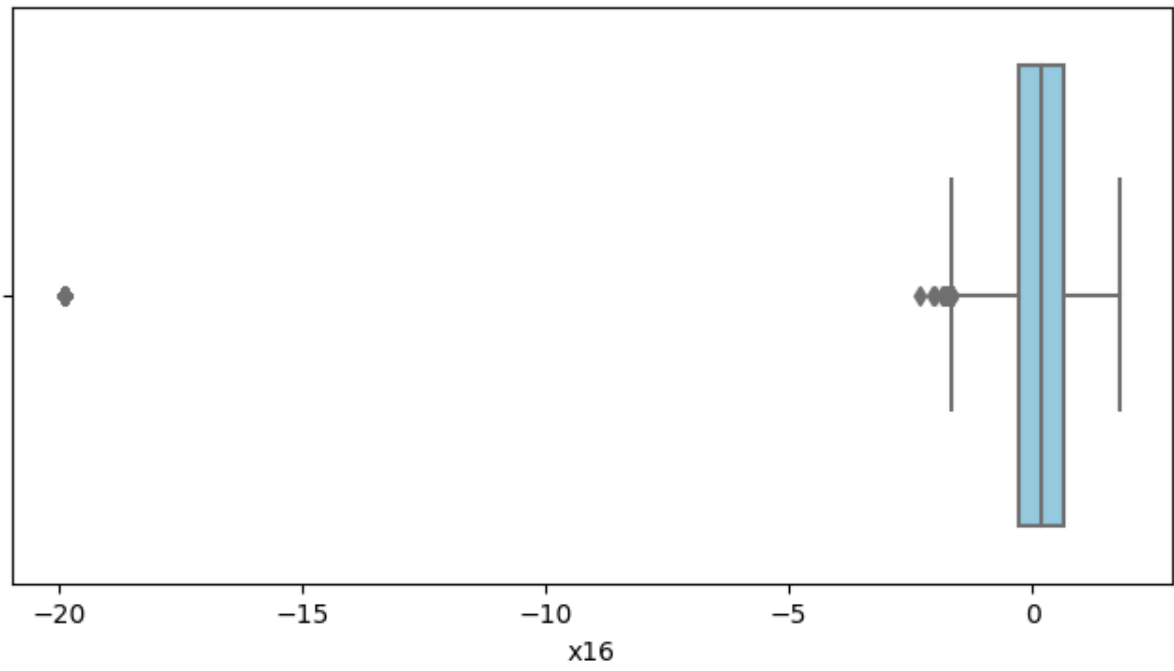
Boxplot for x14



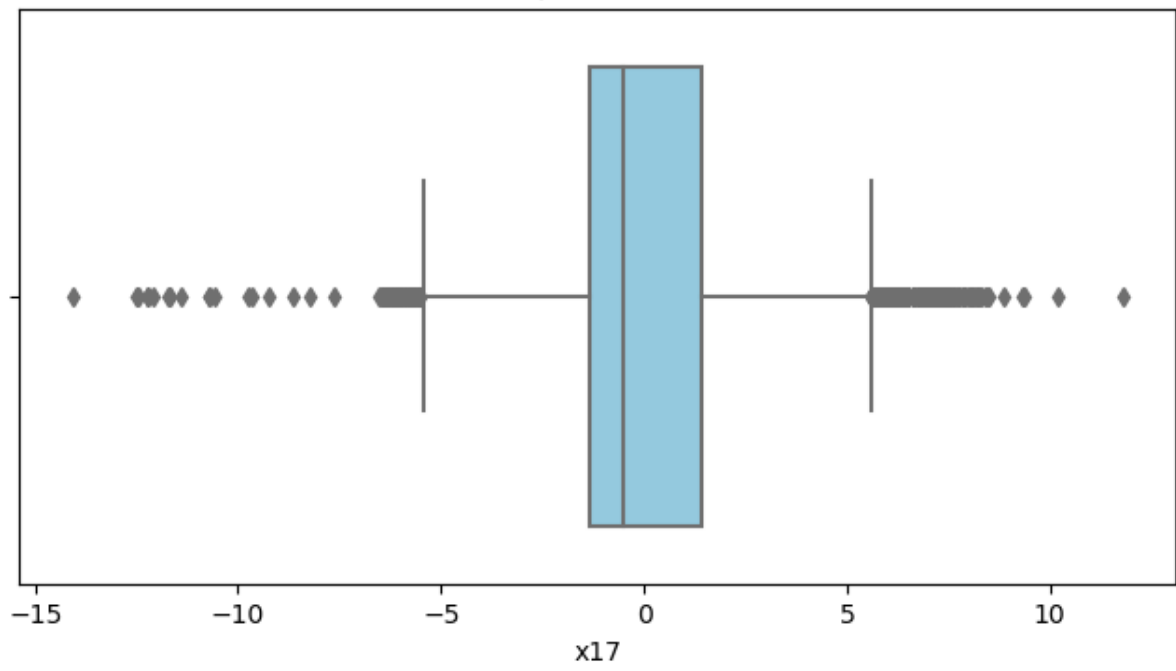
Boxplot for x15



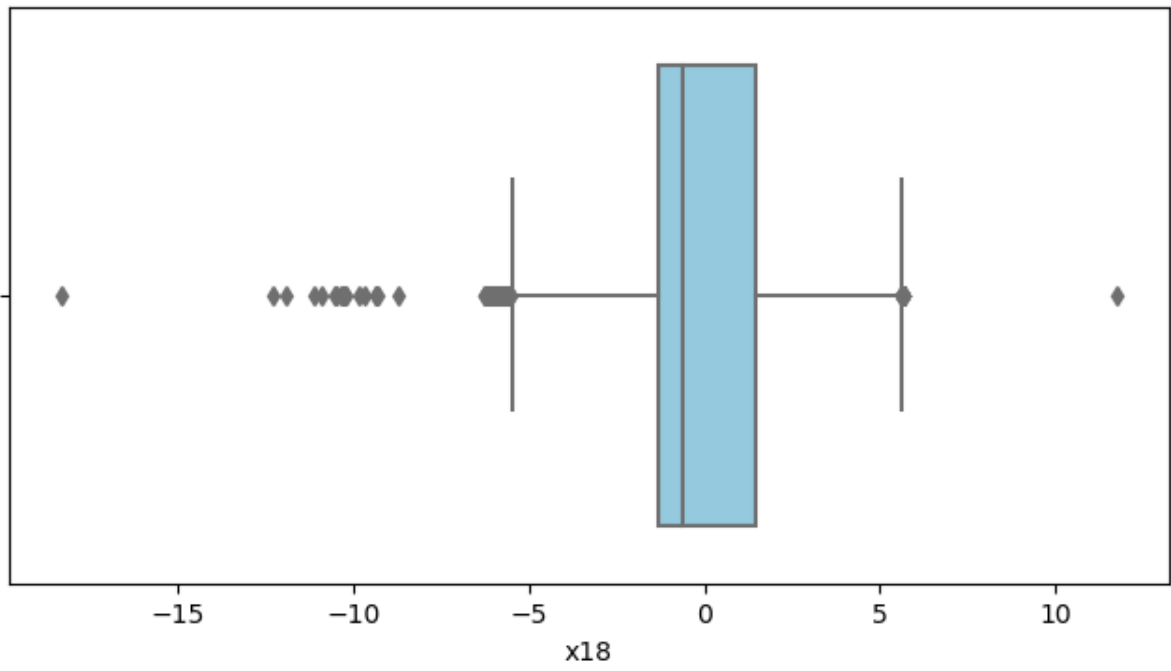
Boxplot for x16



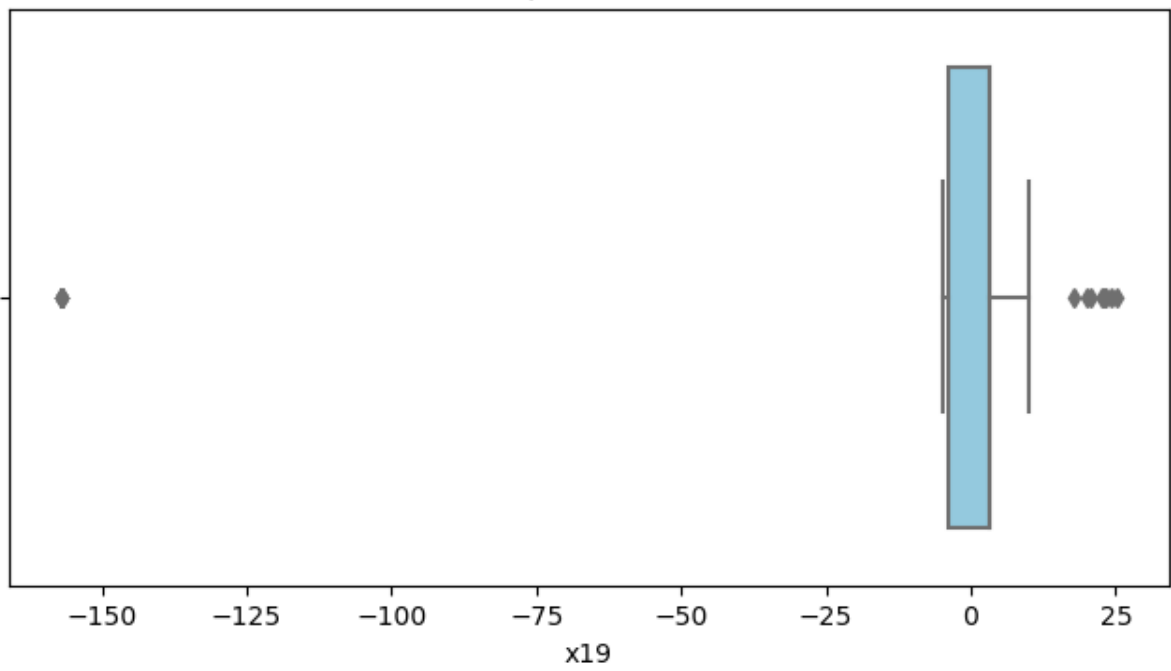
Boxplot for x17



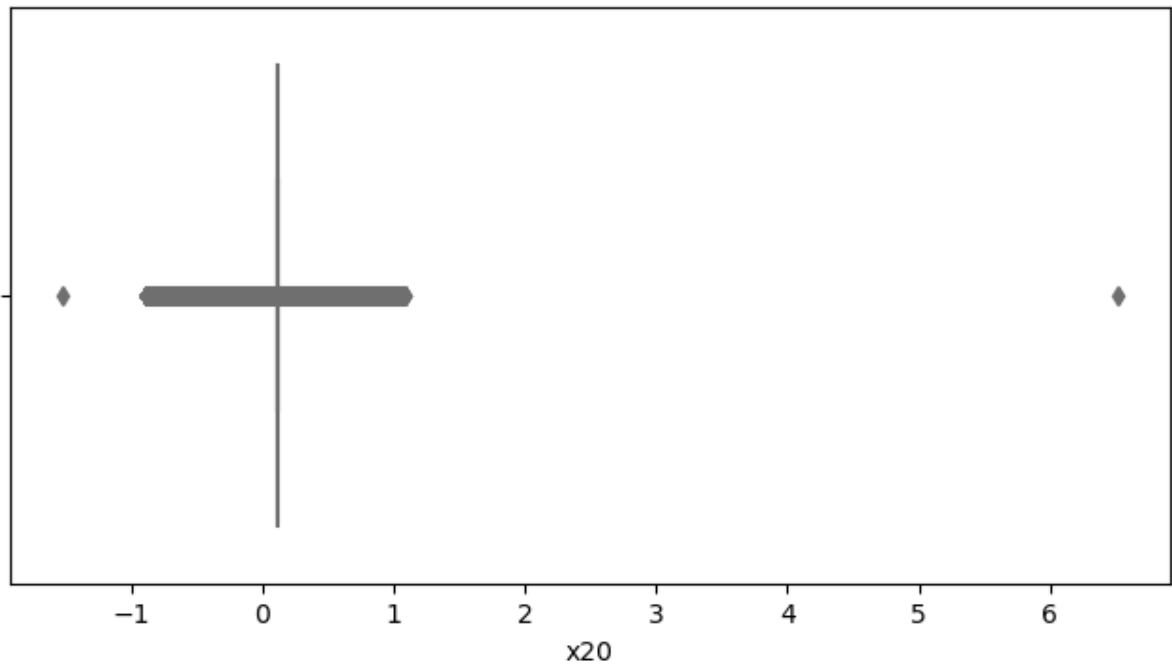
Boxplot for x18



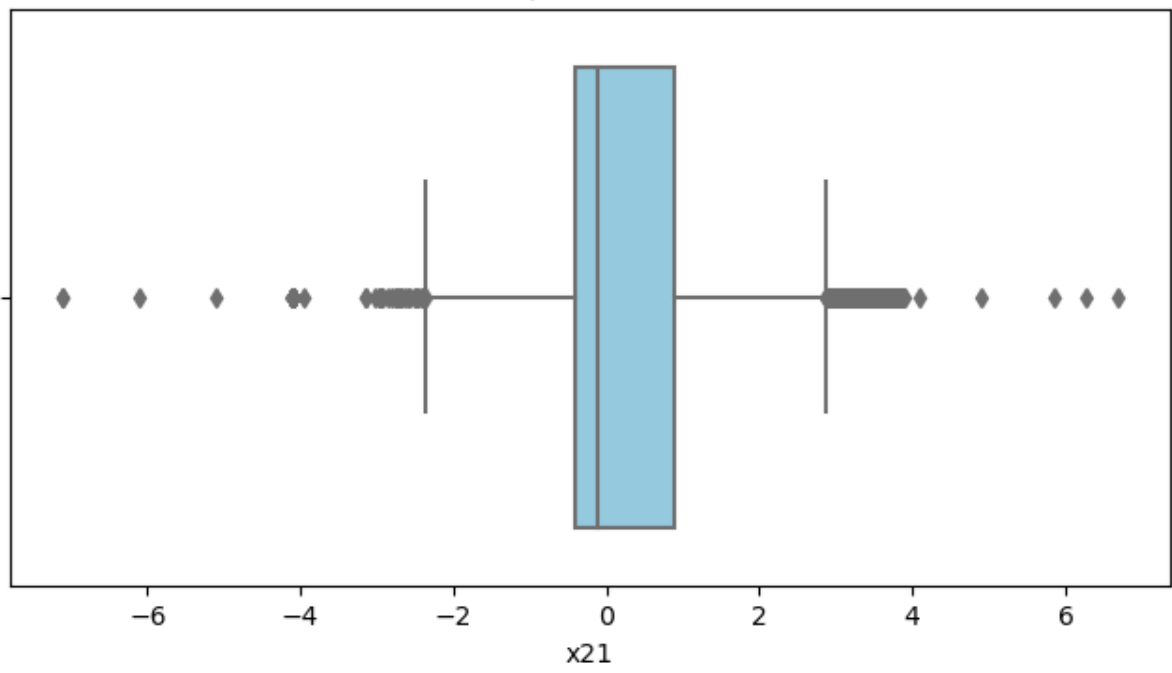
Boxplot for x19



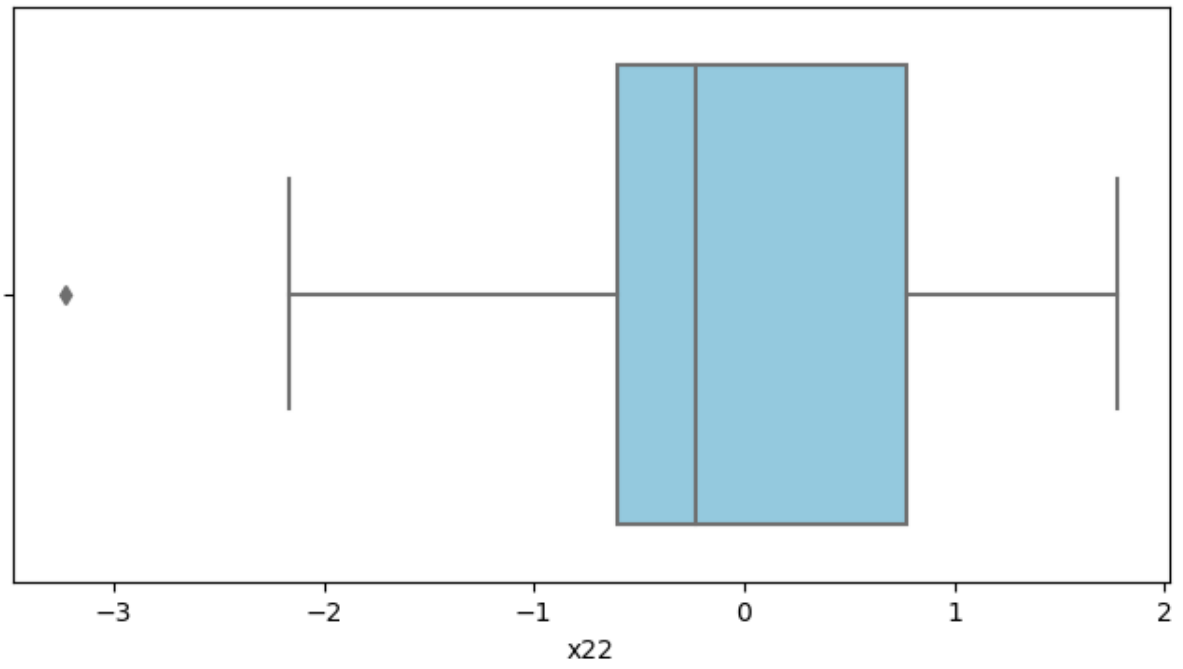
Boxplot for x20



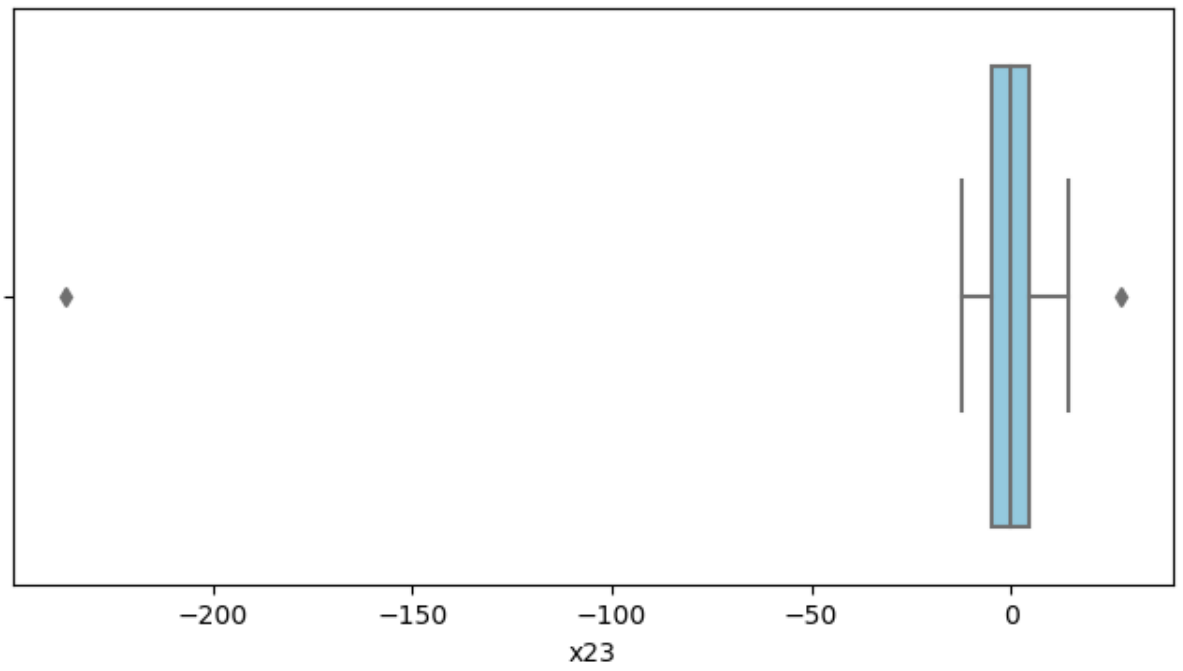
Boxplot for x21



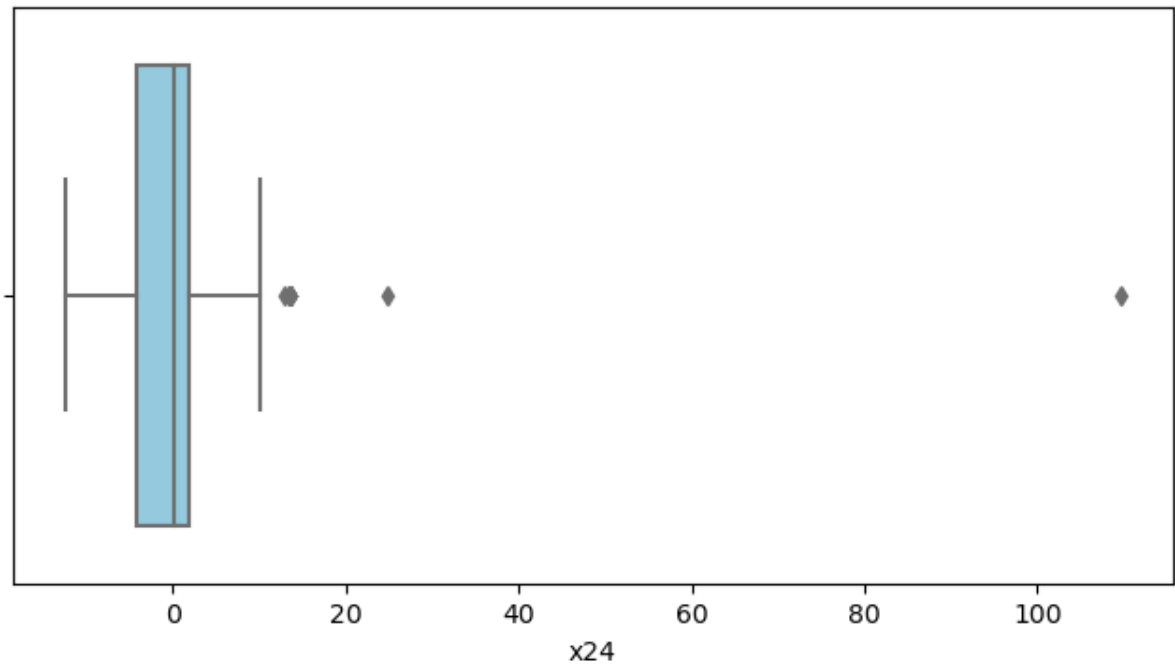
Boxplot for x22



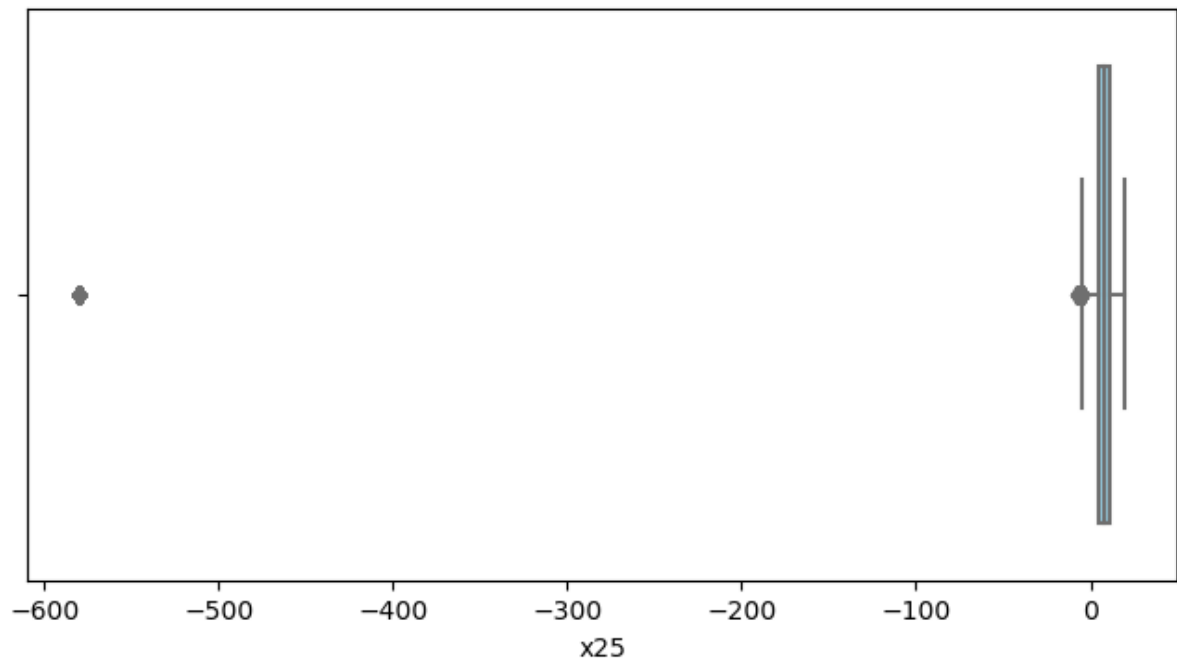
Boxplot for x23



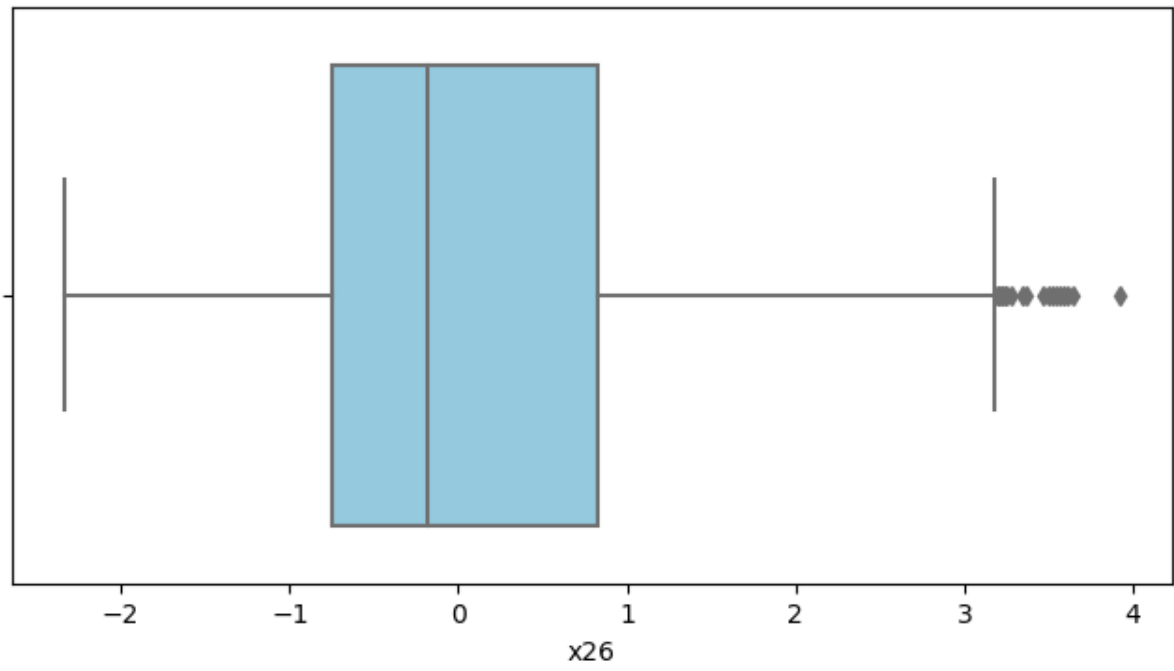
Boxplot for x24



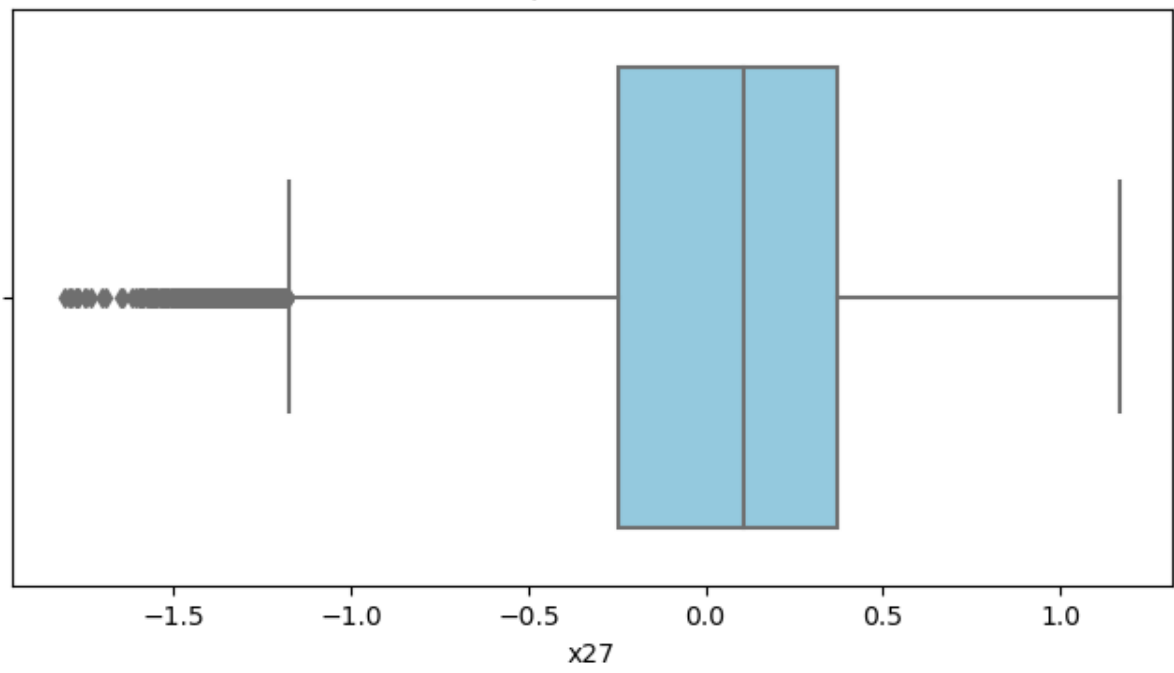
Boxplot for x25



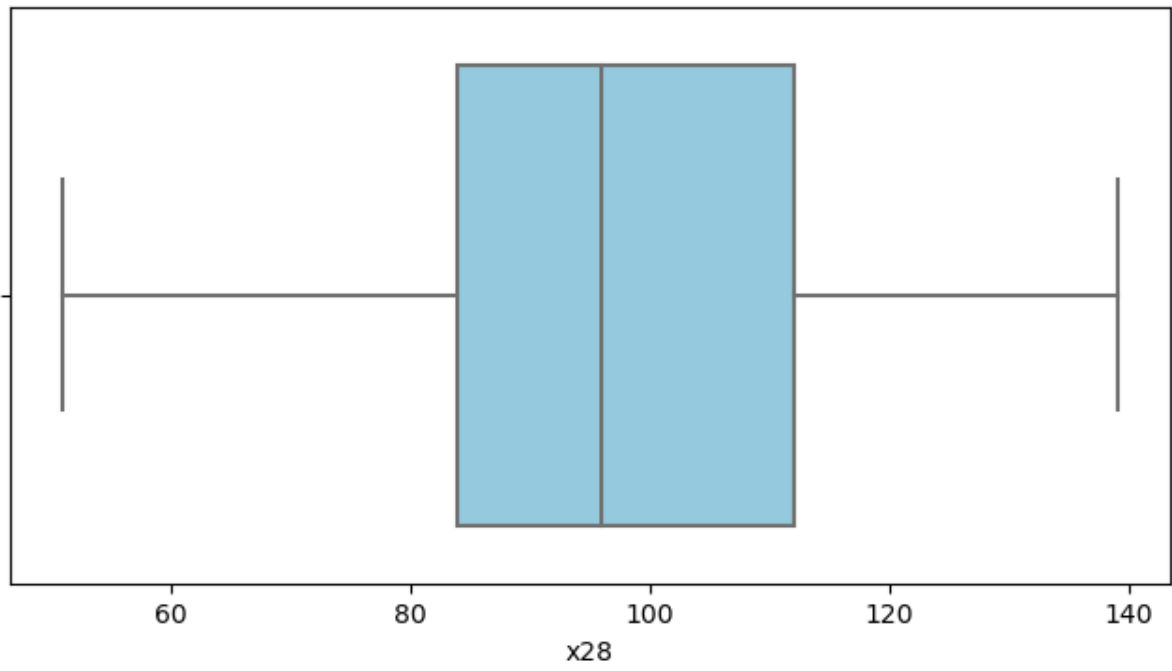
Boxplot for x26



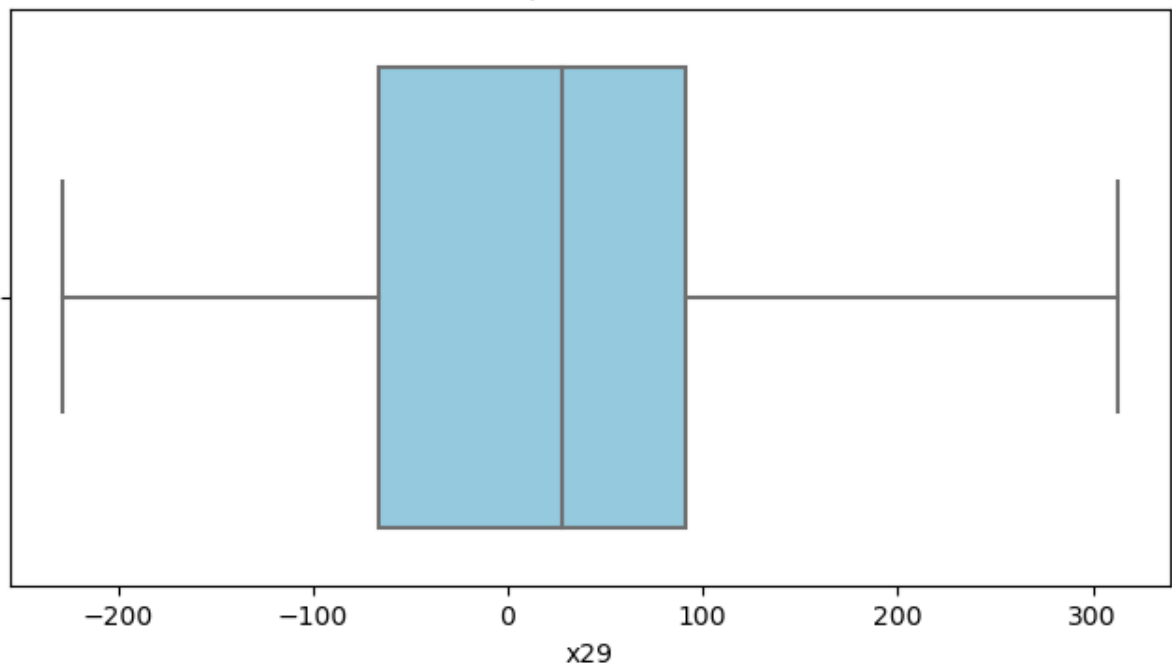
Boxplot for x27



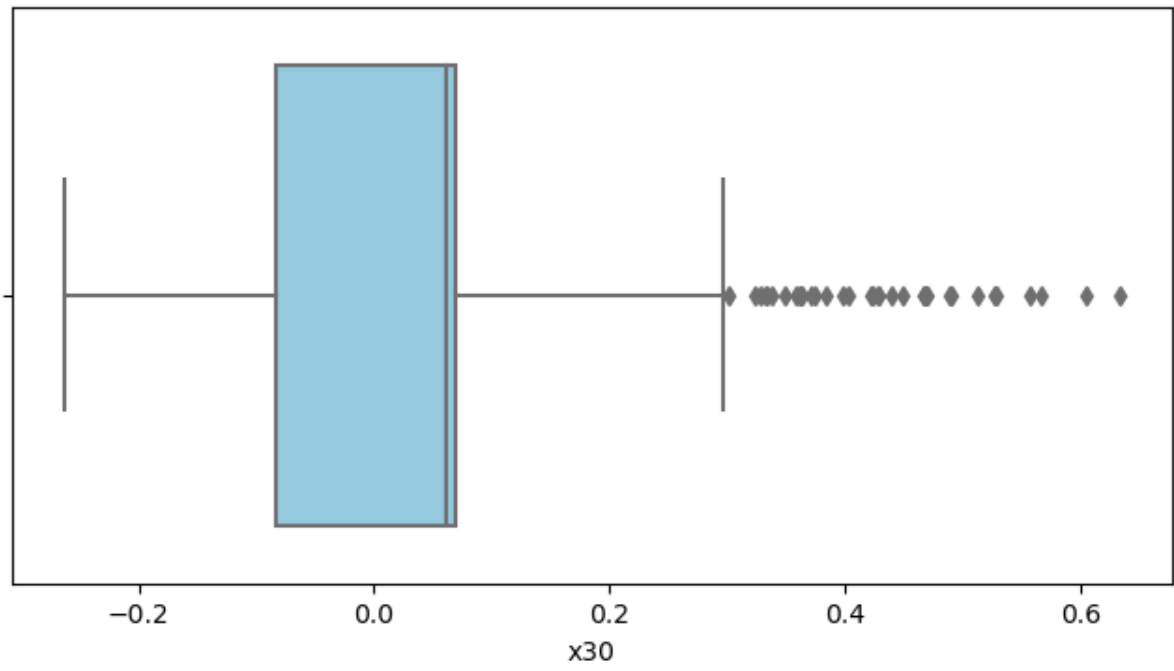
Boxplot for x28



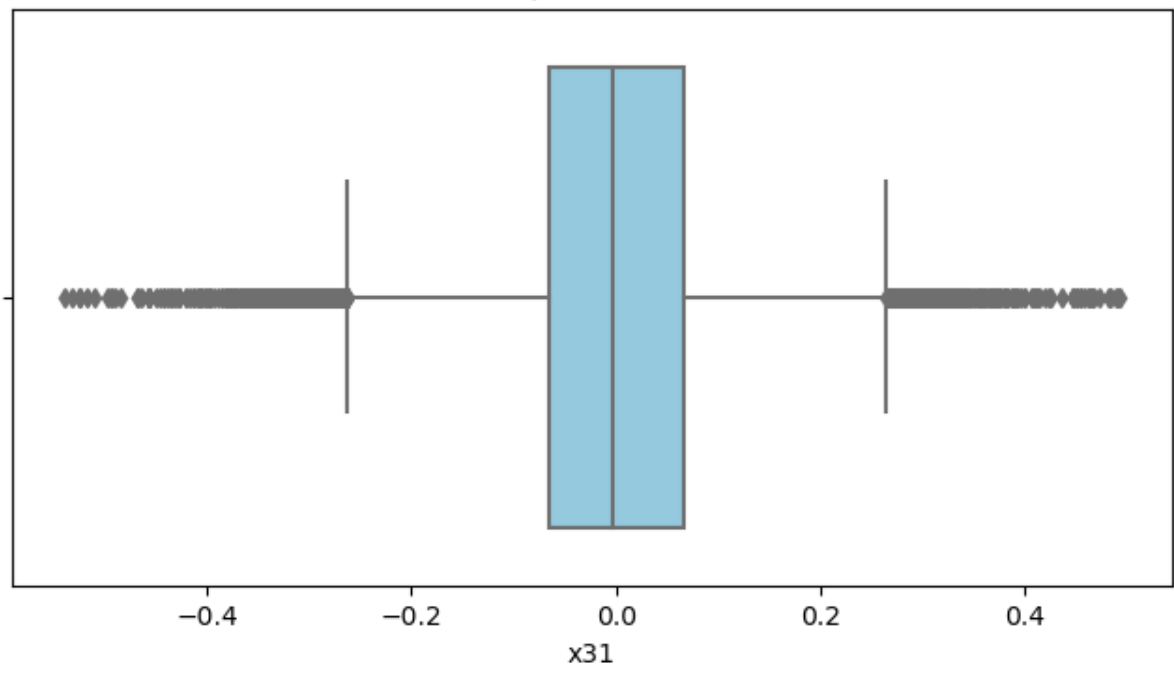
Boxplot for x29



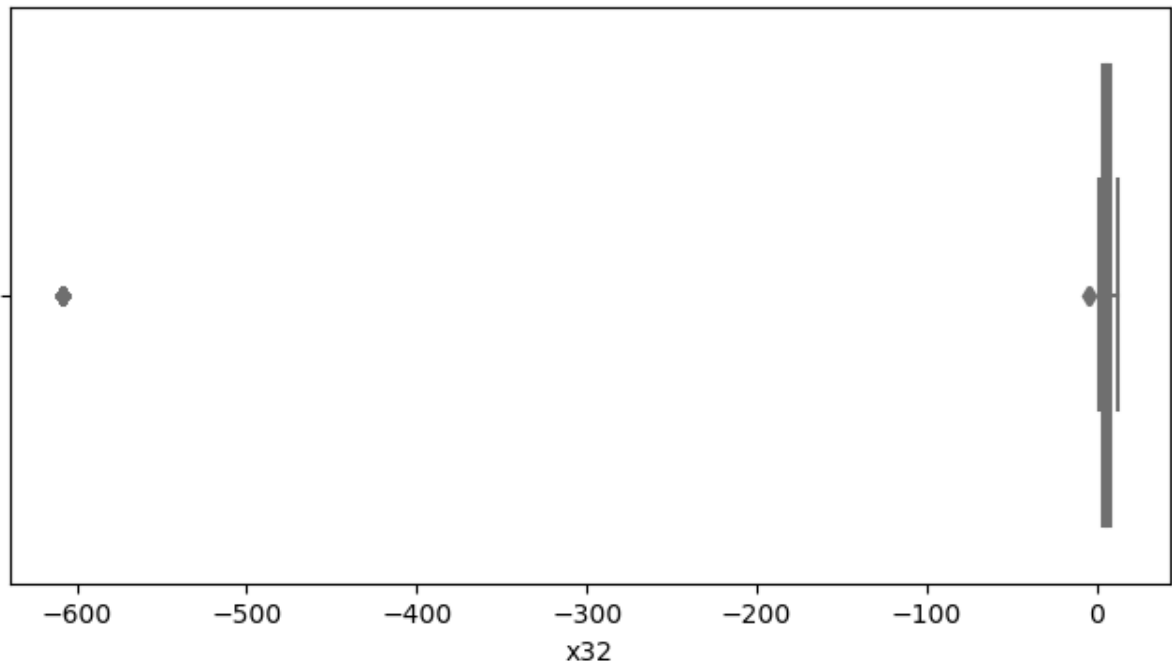
Boxplot for x30



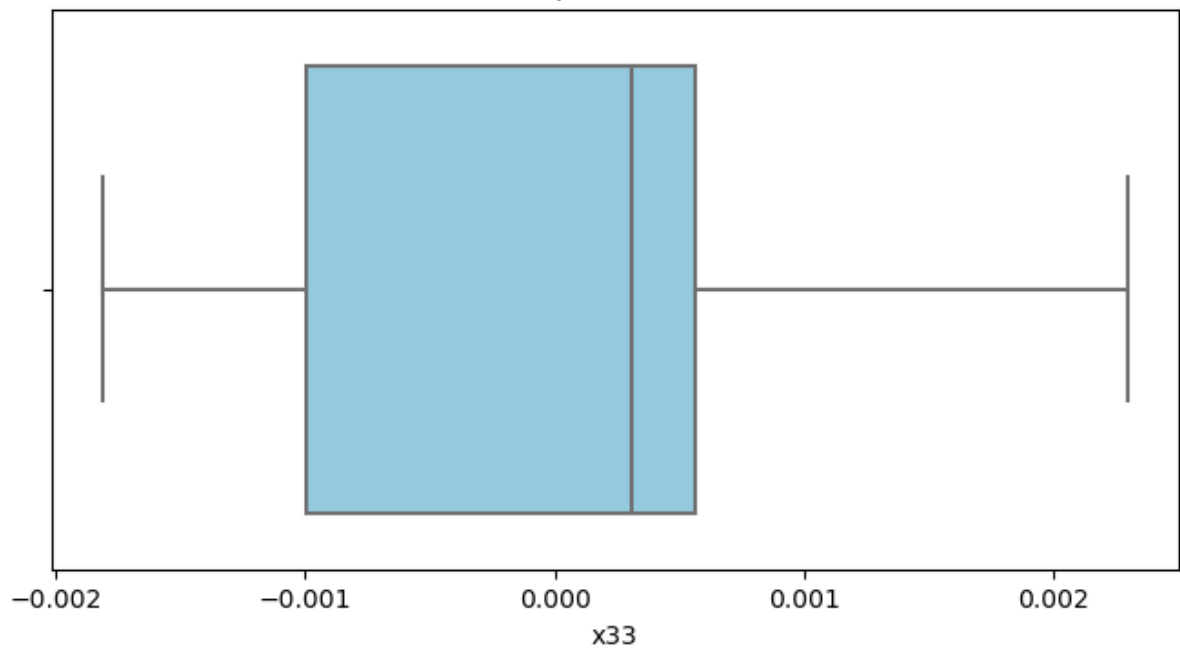
Boxplot for x31



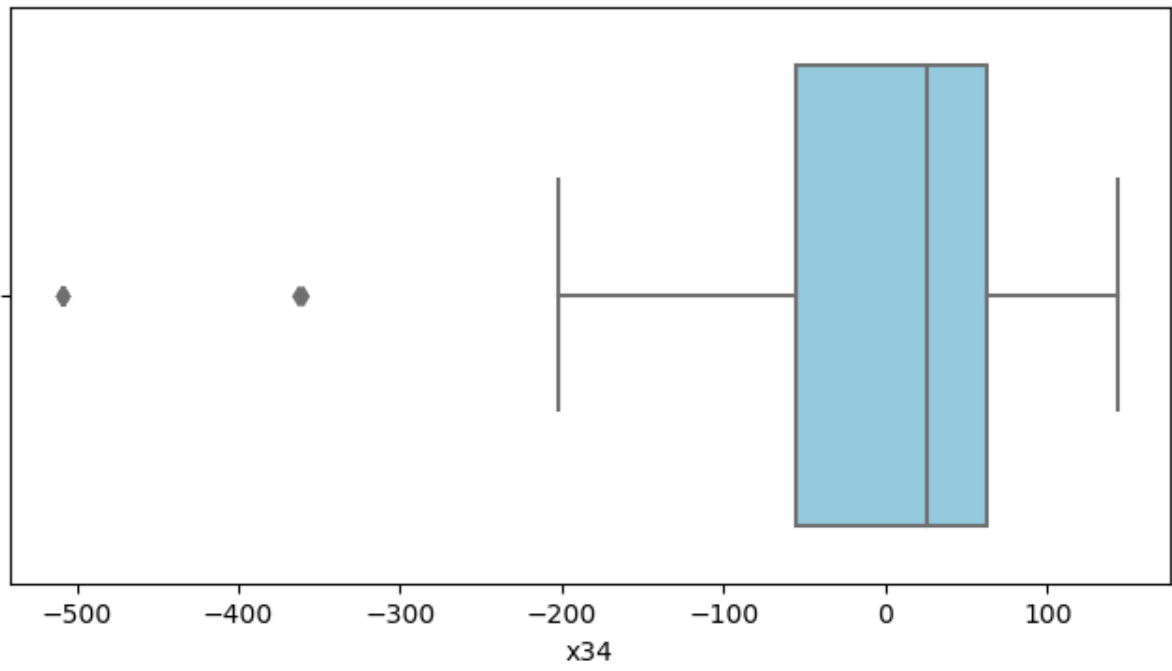
Boxplot for x32



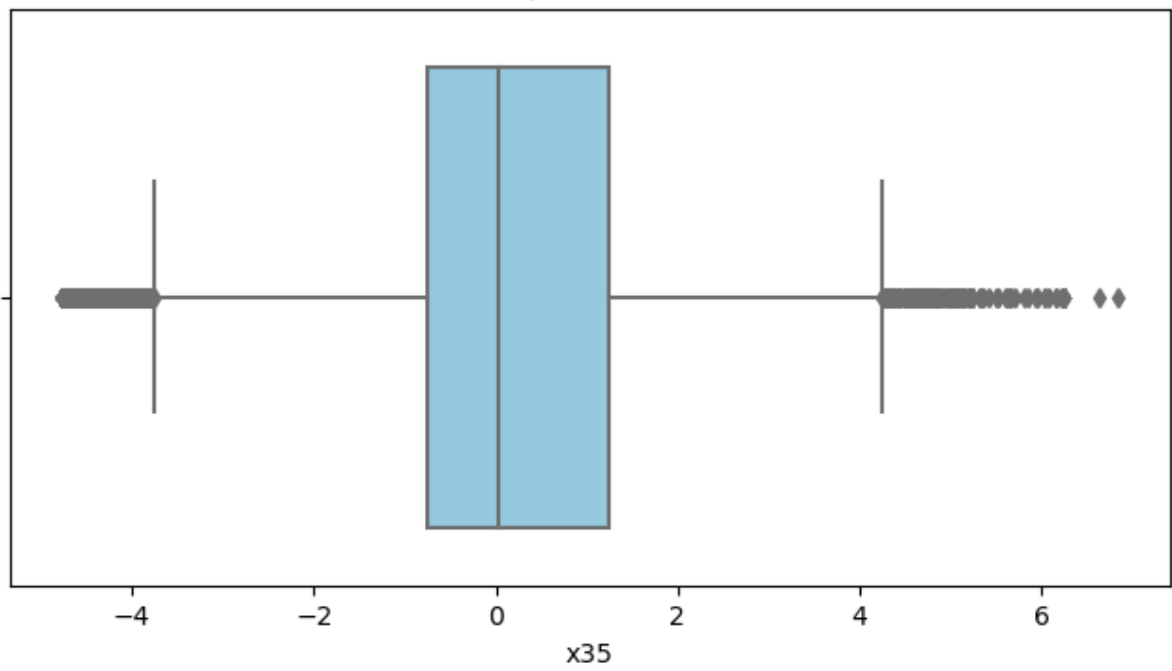
Boxplot for x33



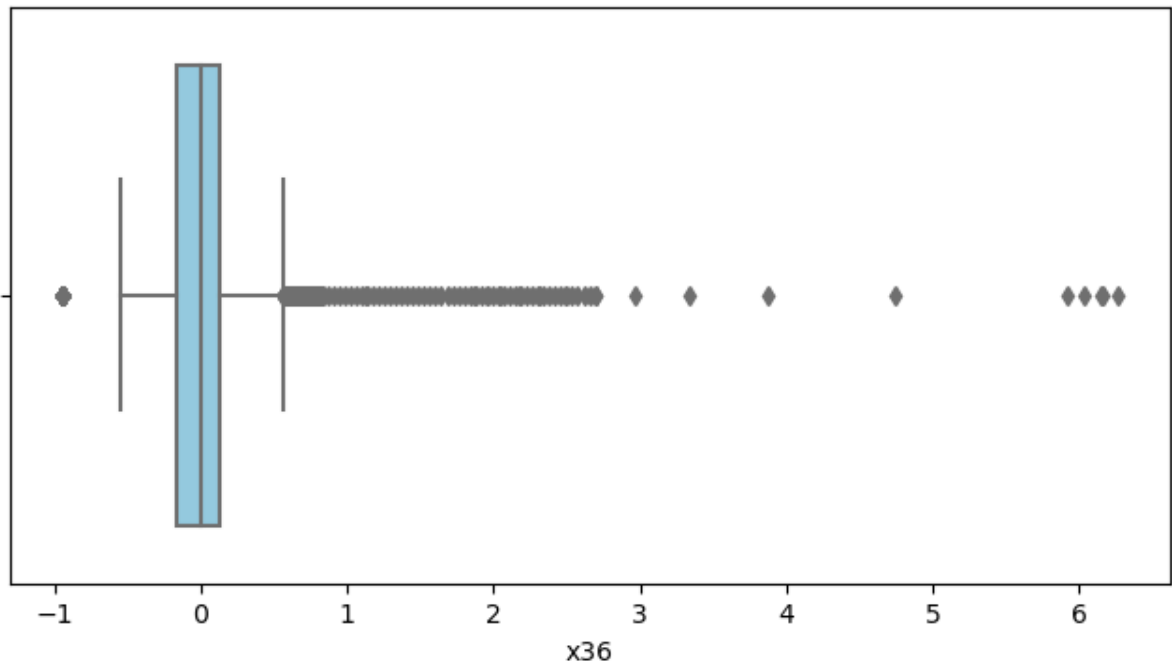
Boxplot for x34



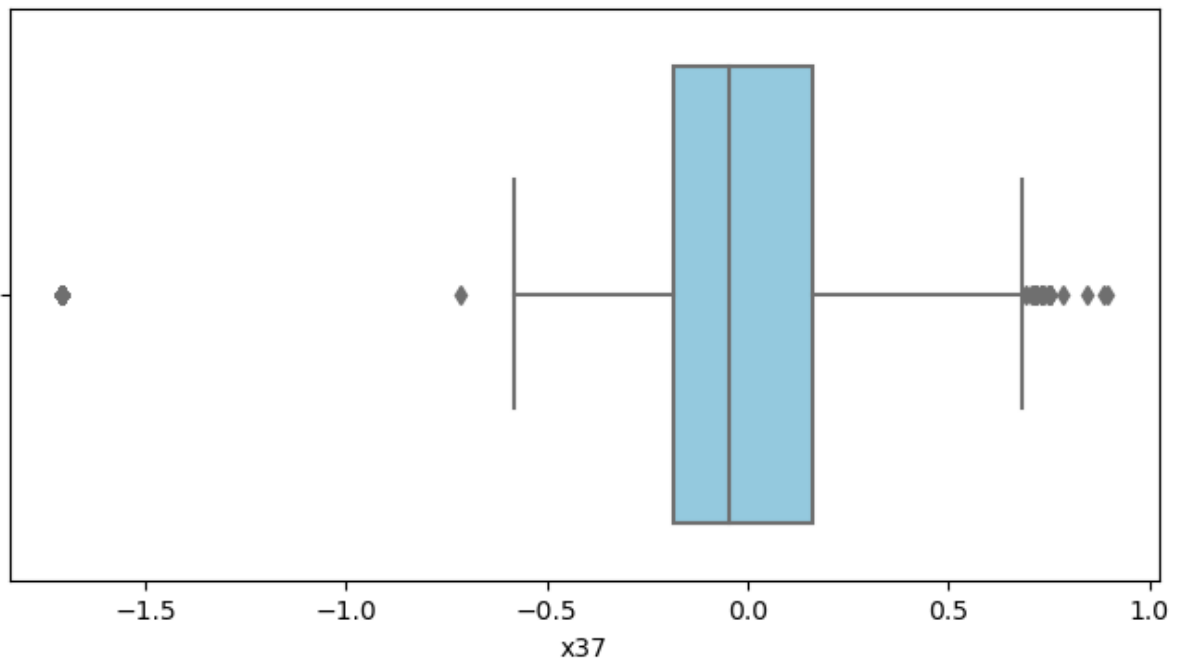
Boxplot for x35



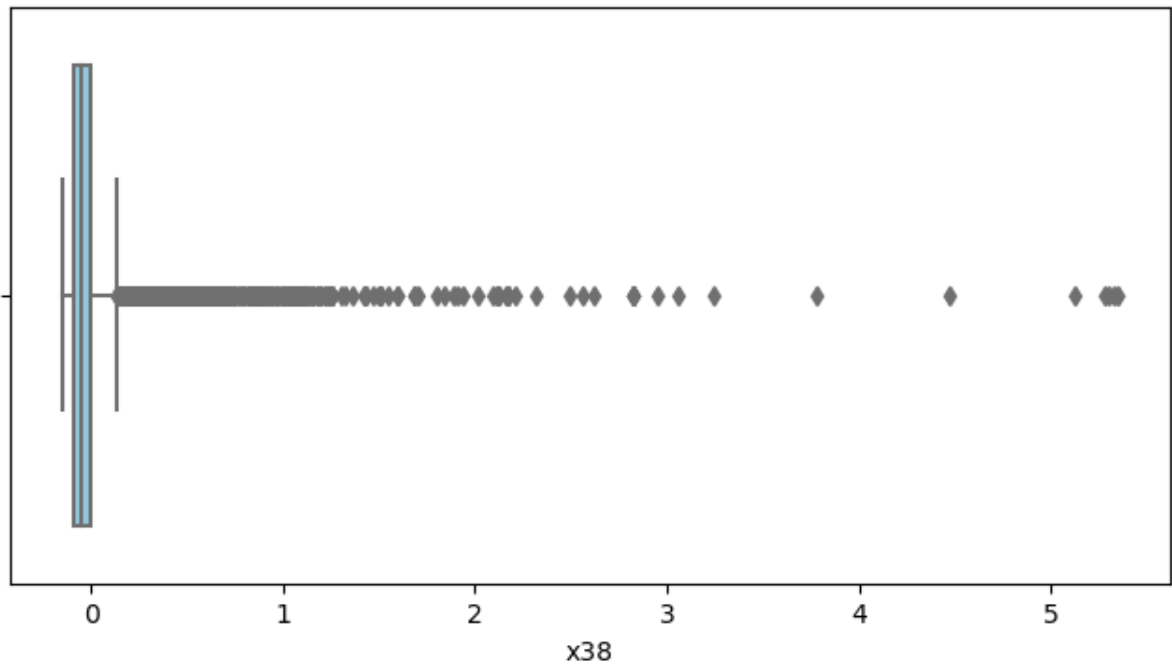
Boxplot for x36



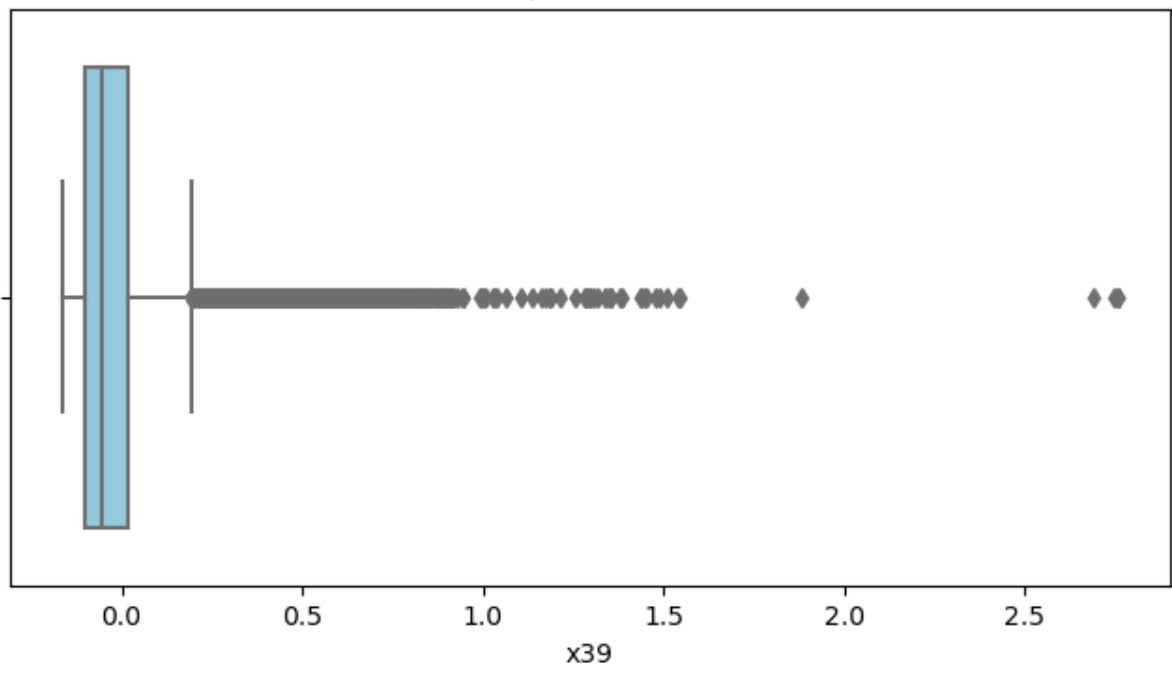
Boxplot for x37



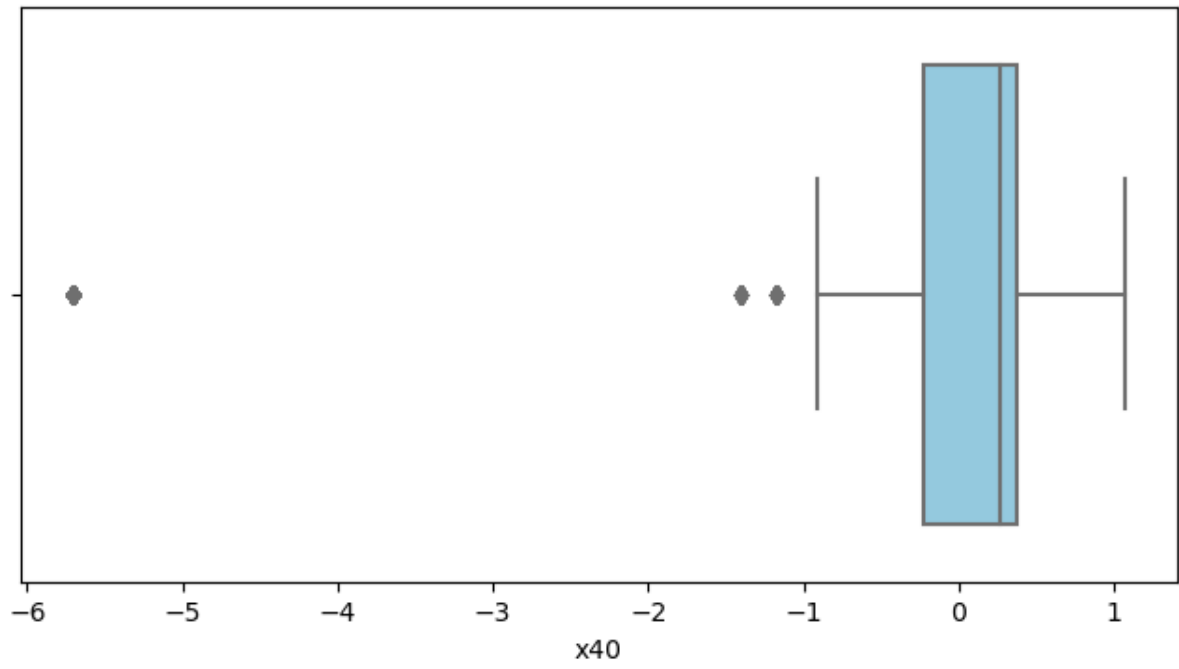
Boxplot for x38



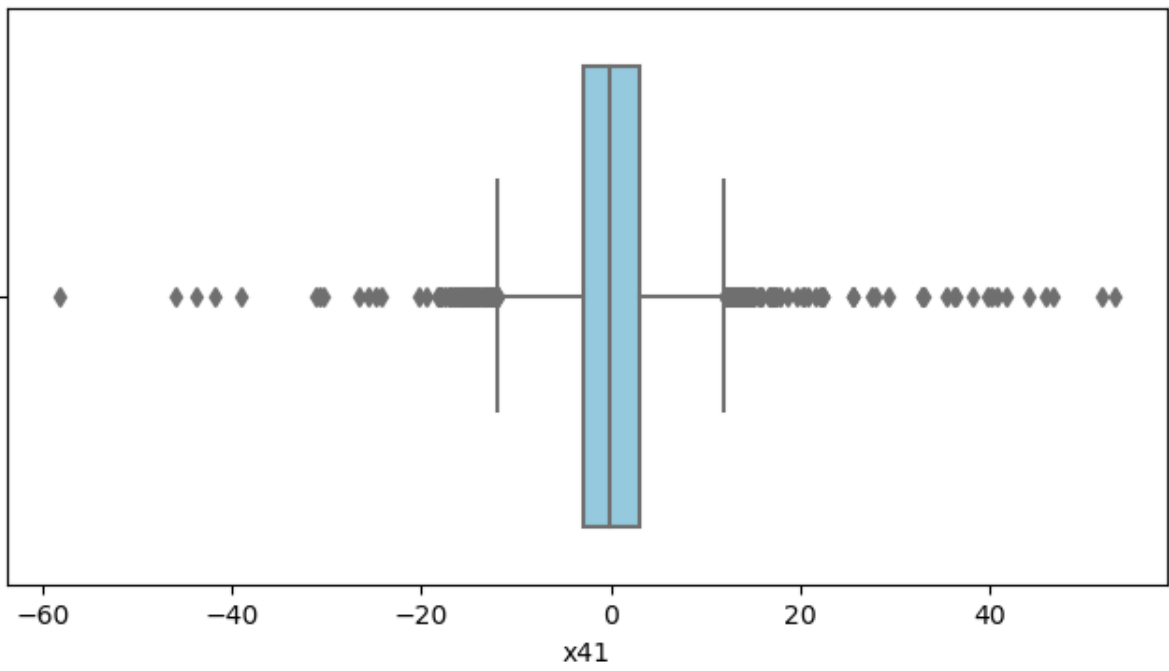
Boxplot for x39



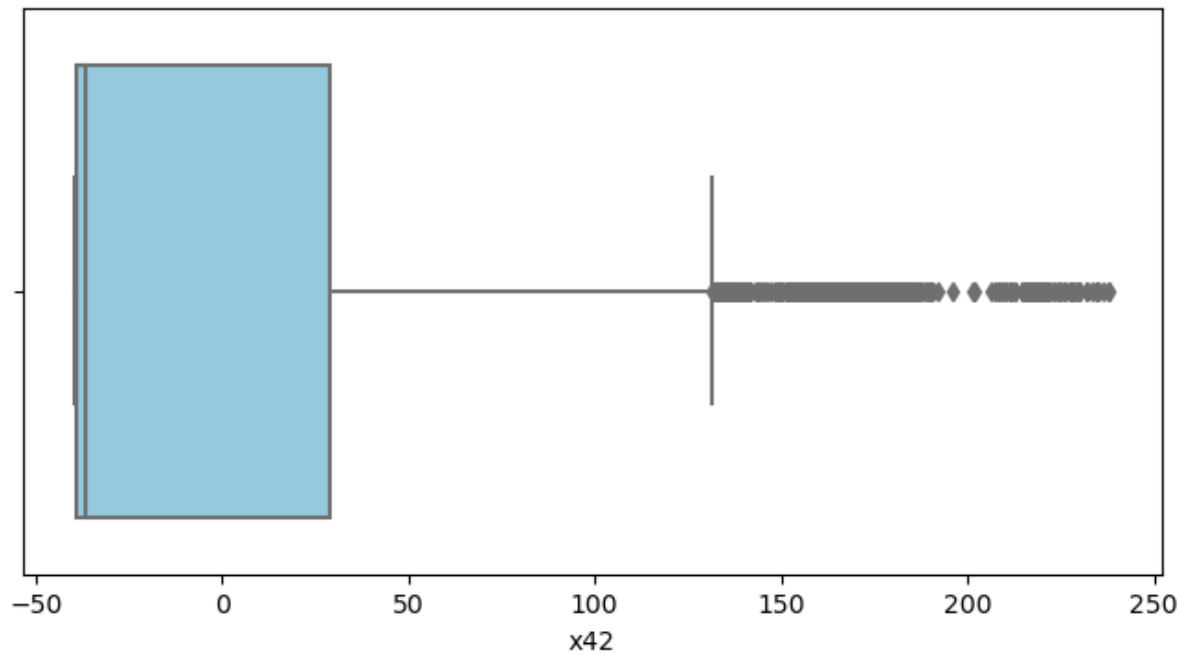
Boxplot for x40



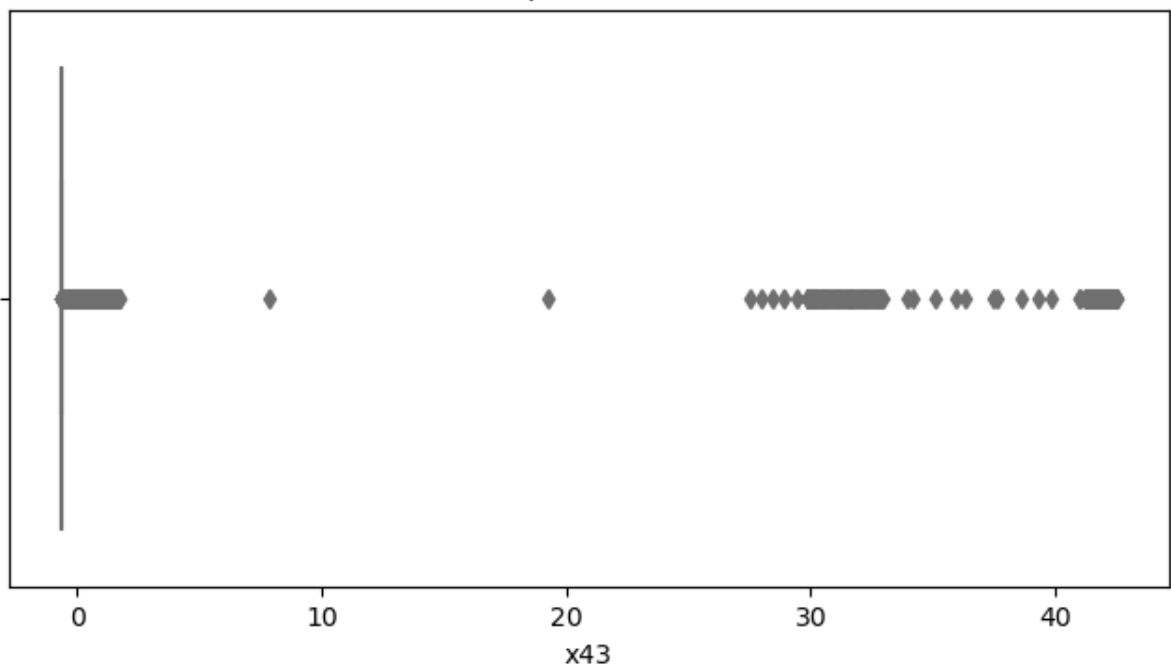
Boxplot for x41



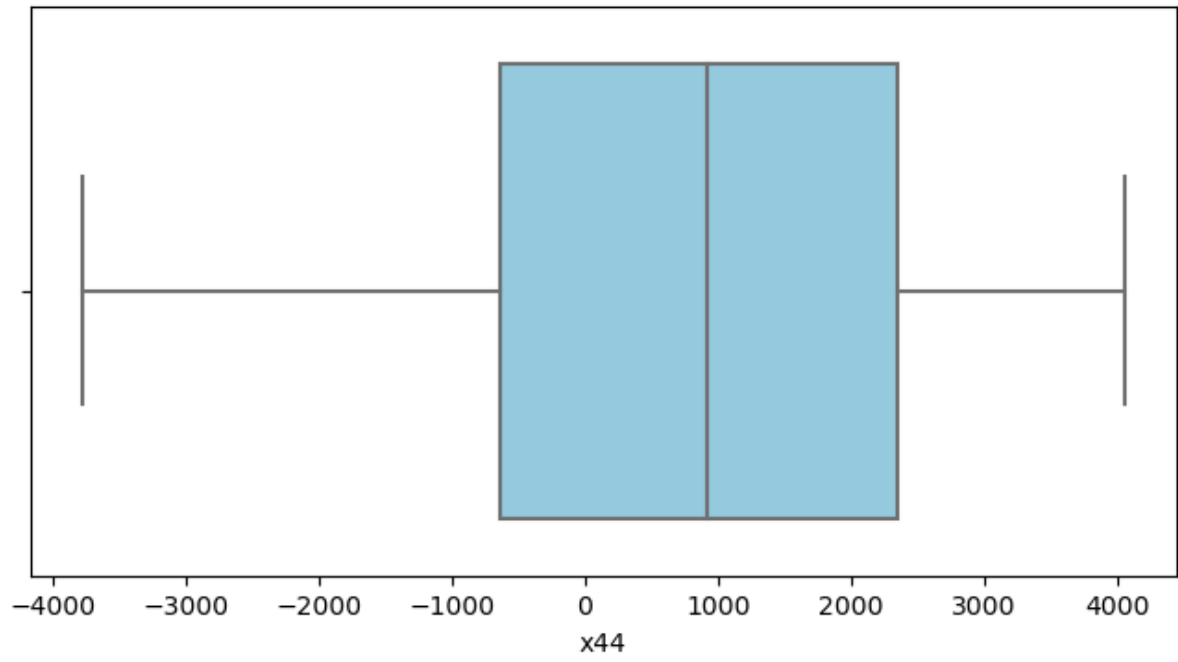
Boxplot for x42



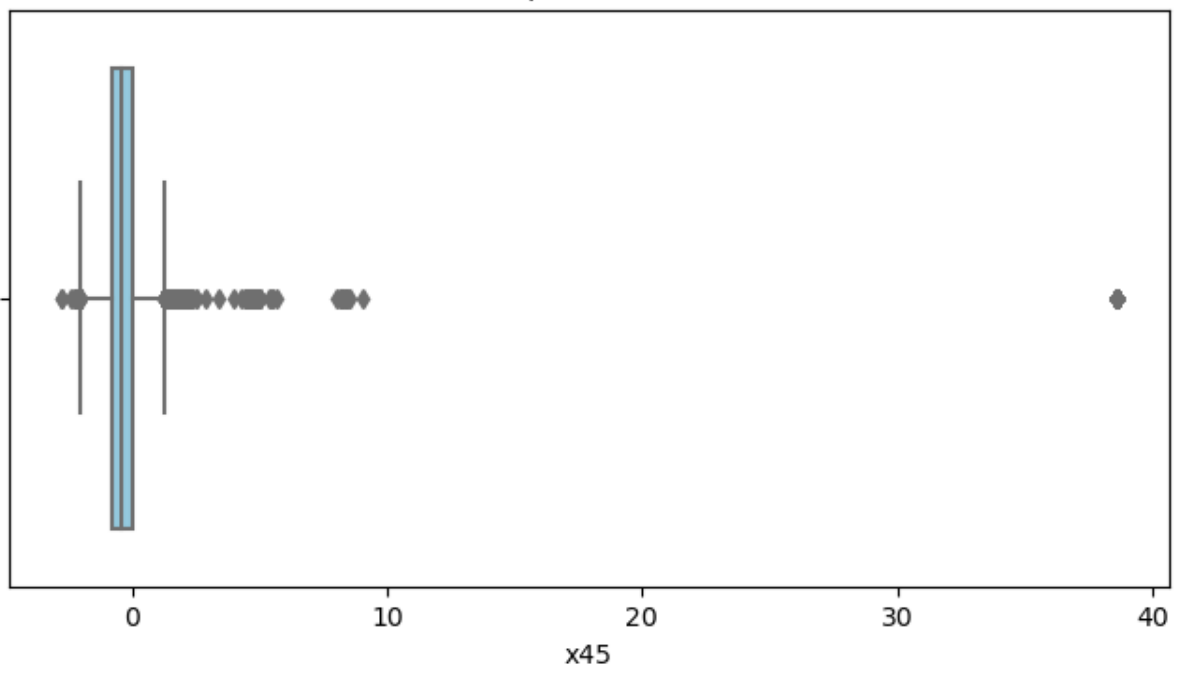
Boxplot for x43



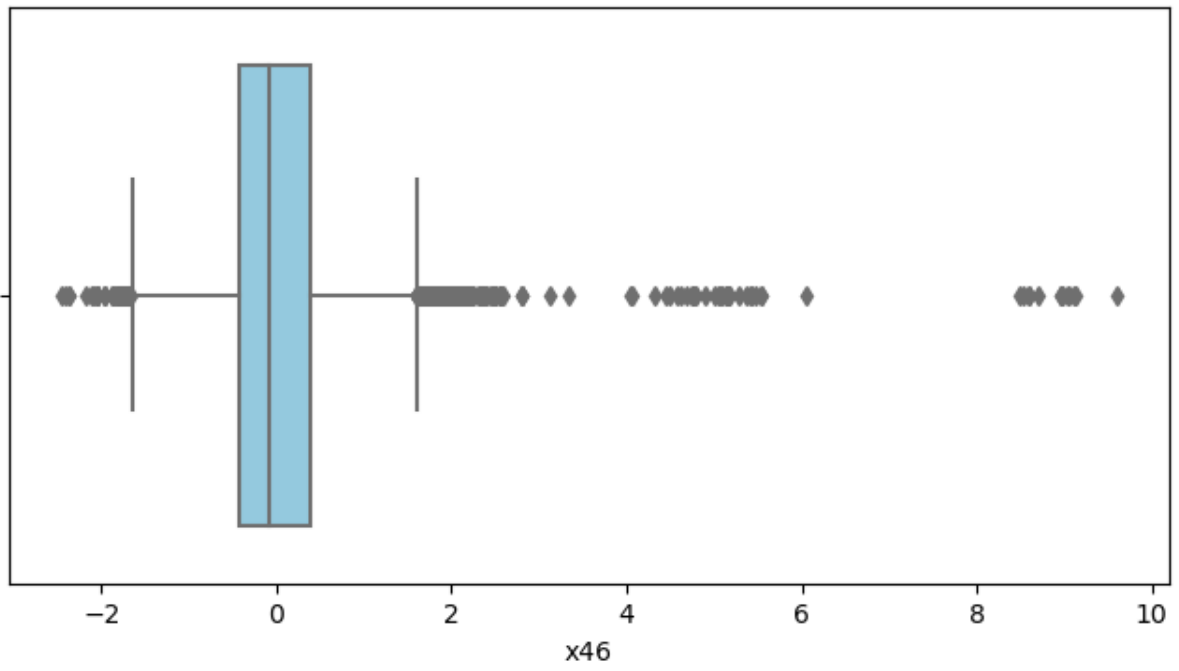
Boxplot for x44



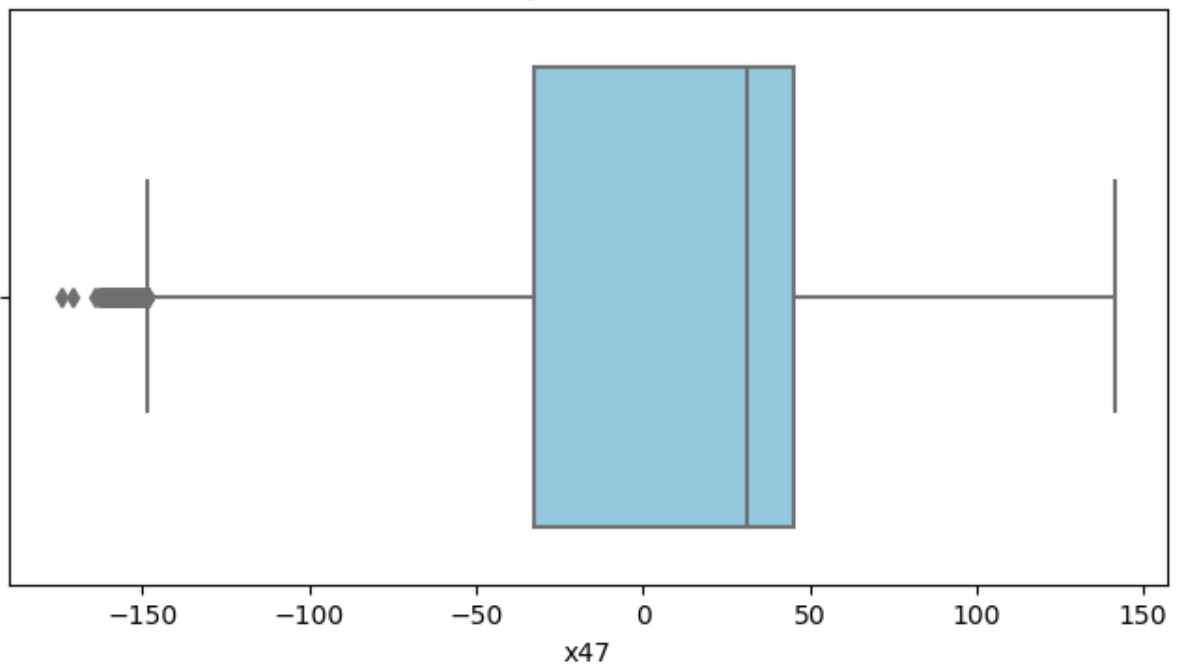
Boxplot for x45



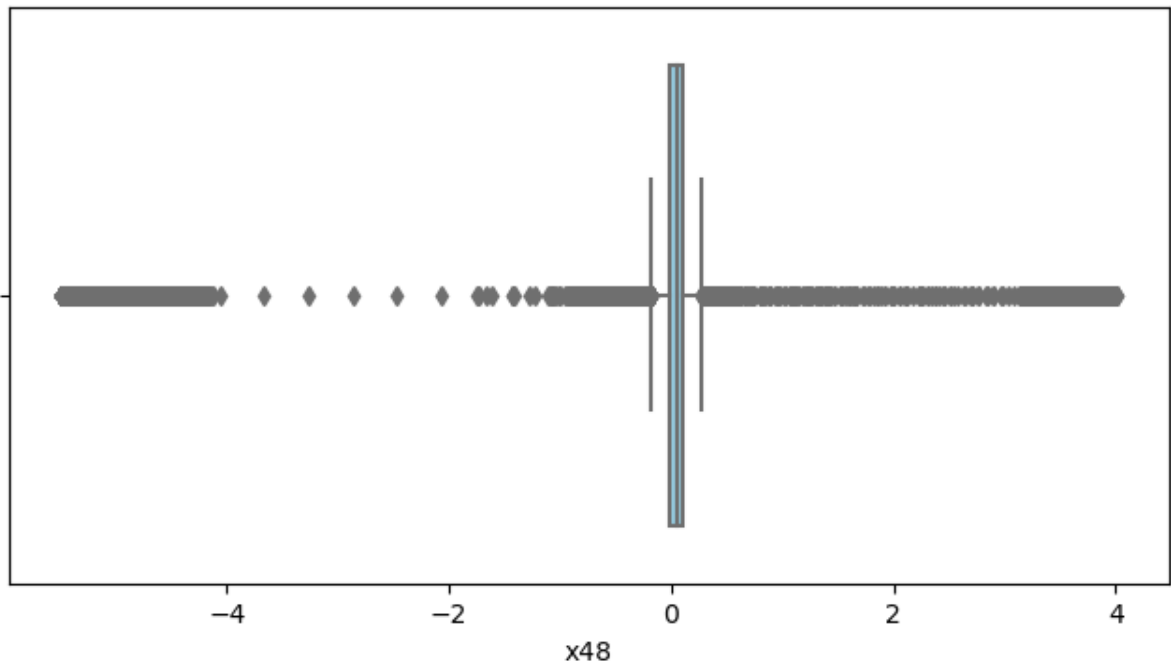
Boxplot for x46



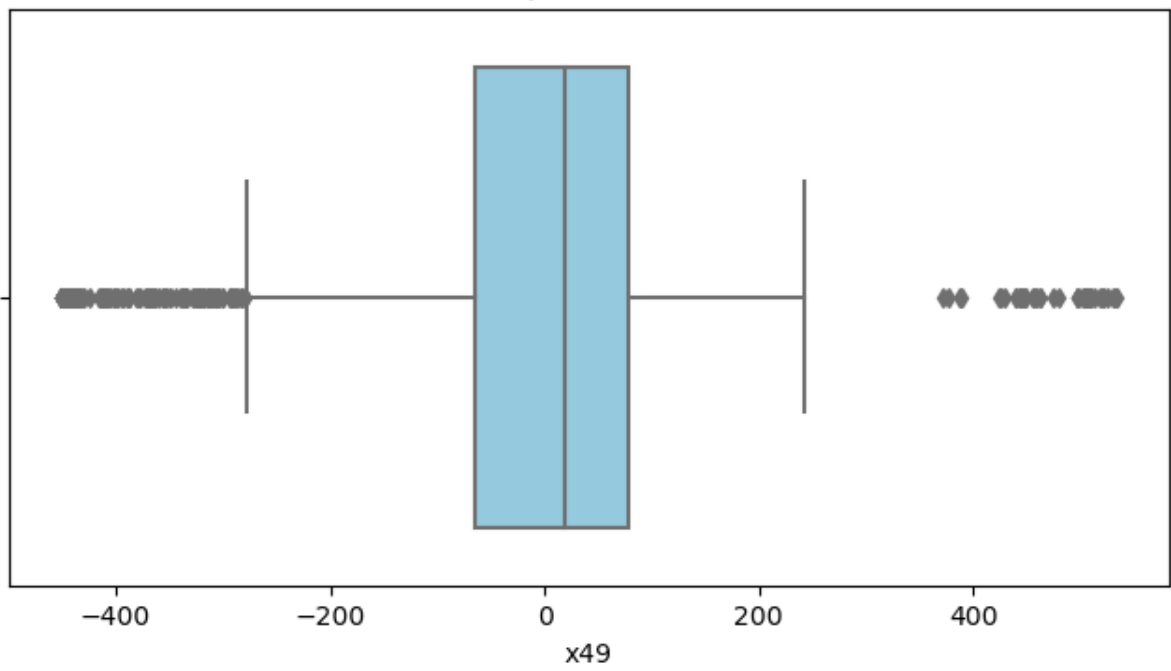
Boxplot for x47



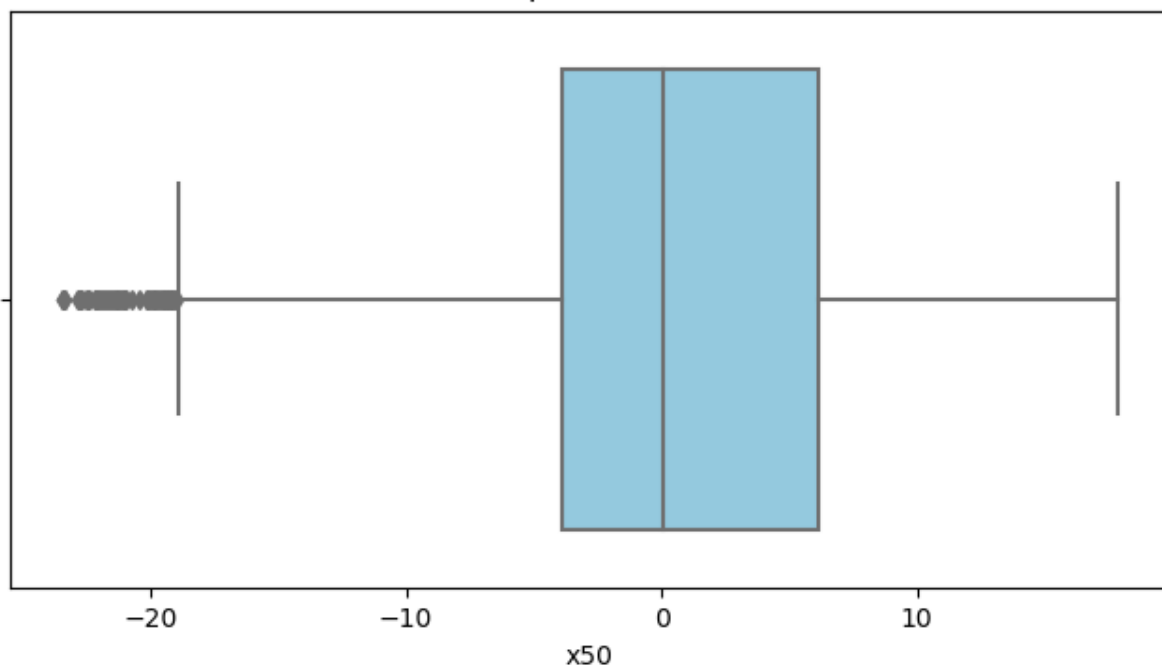
Boxplot for x48



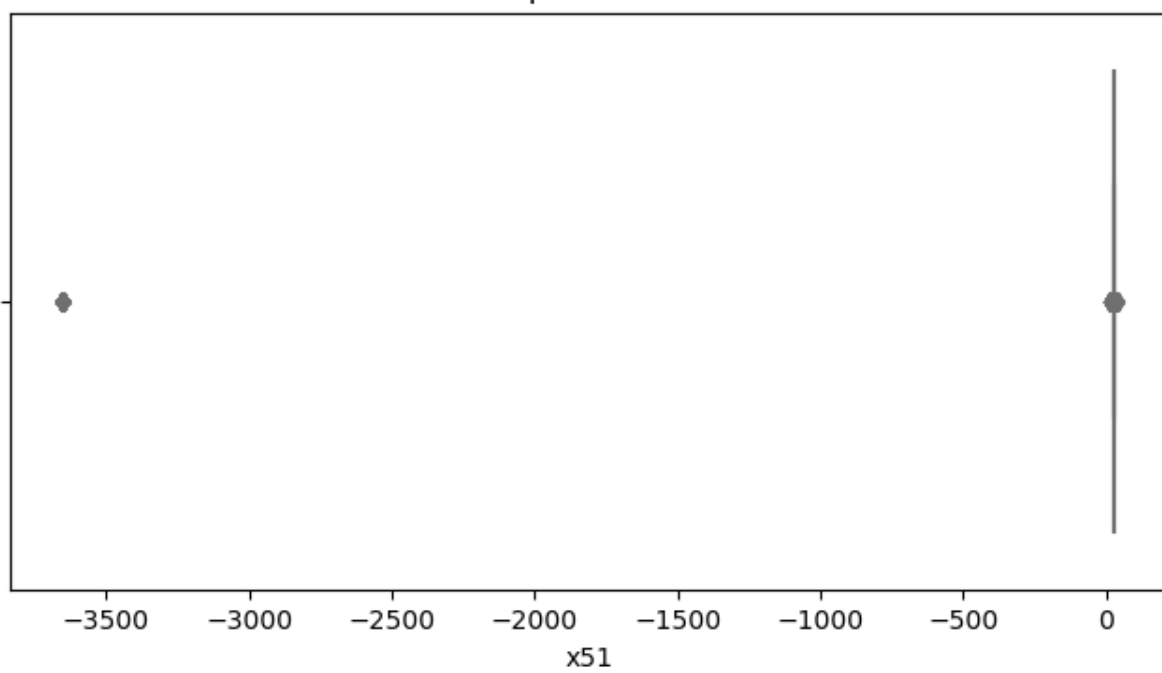
Boxplot for x49



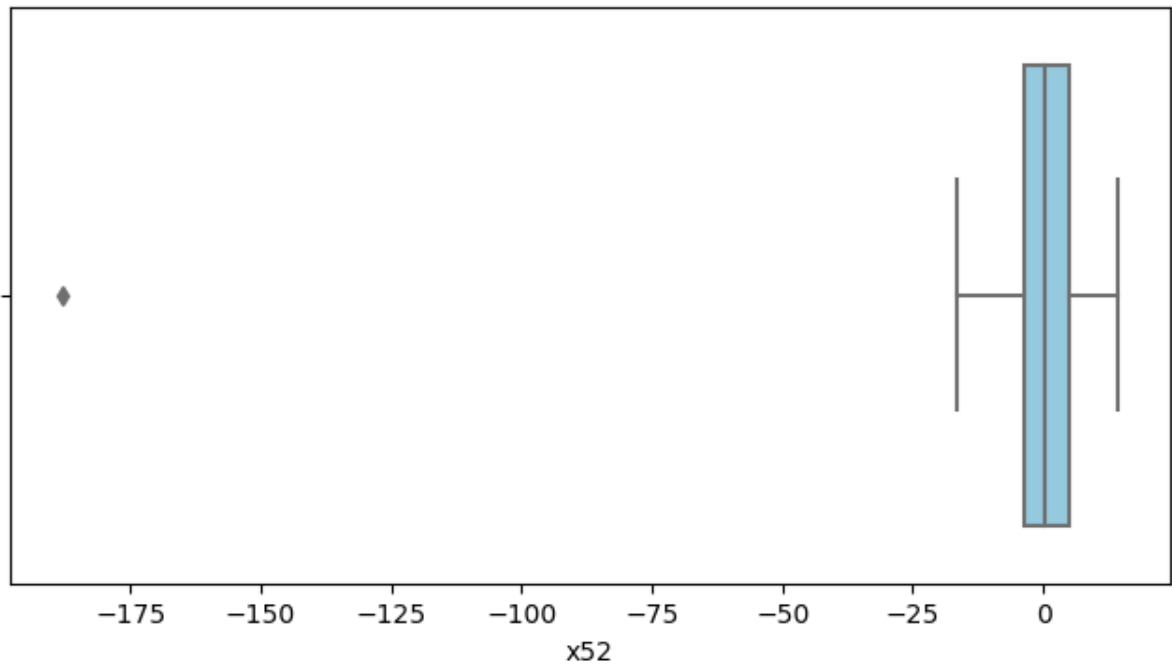
Boxplot for x50



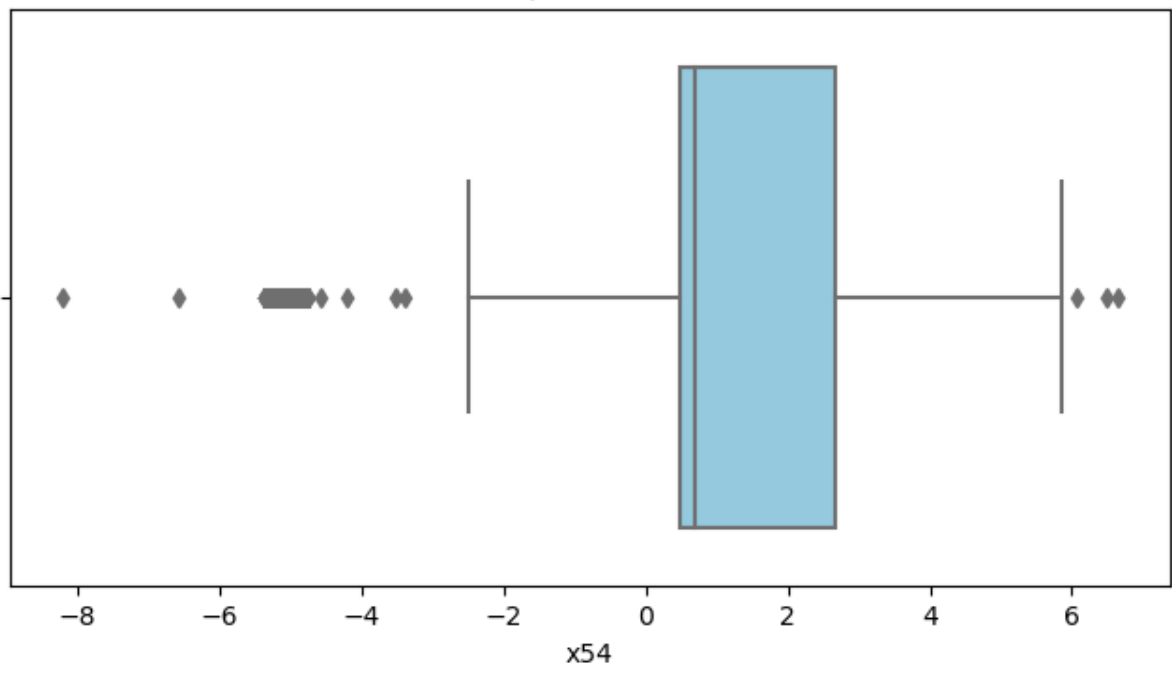
Boxplot for x51



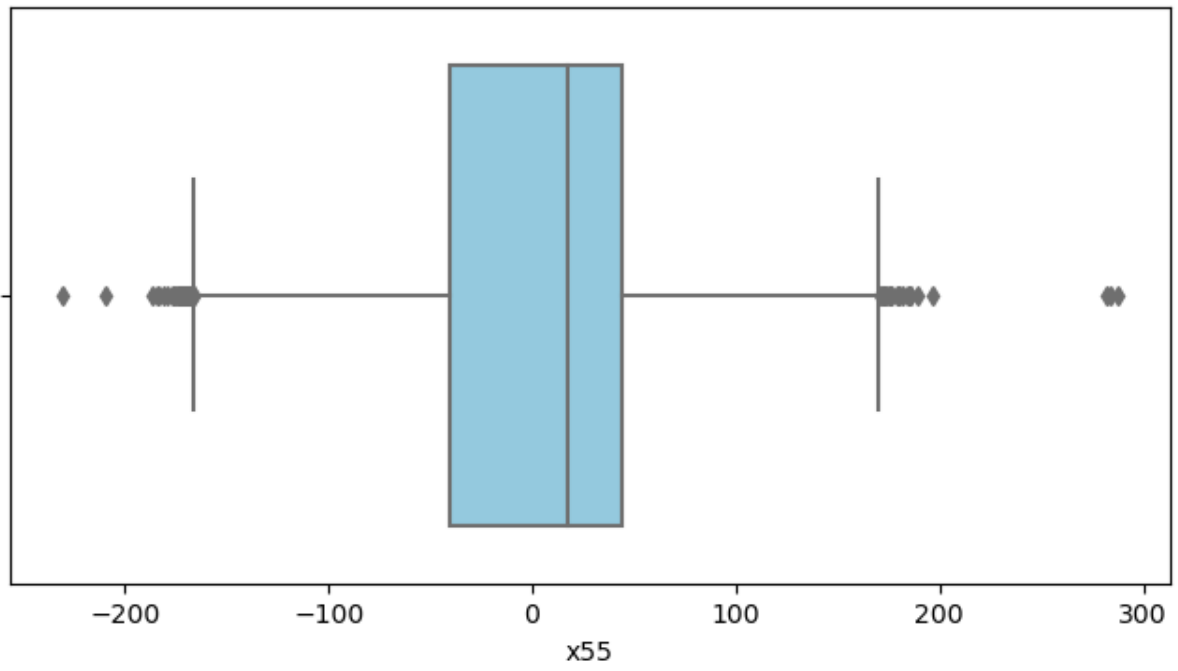
Boxplot for x52



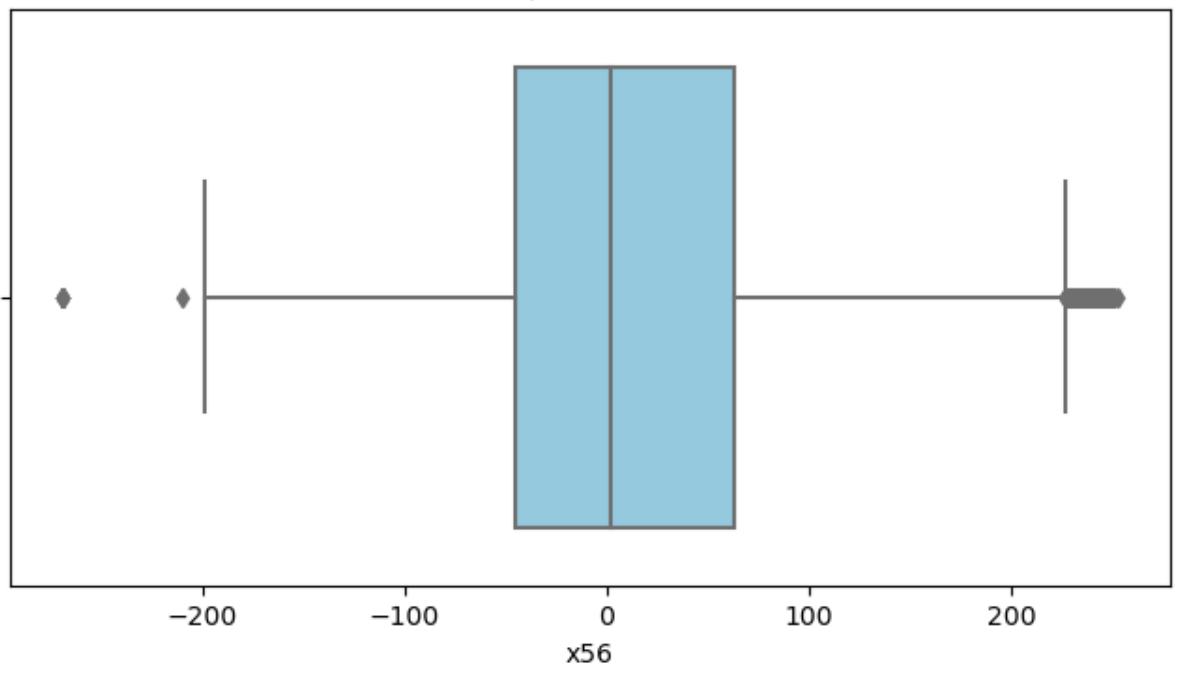
Boxplot for x54



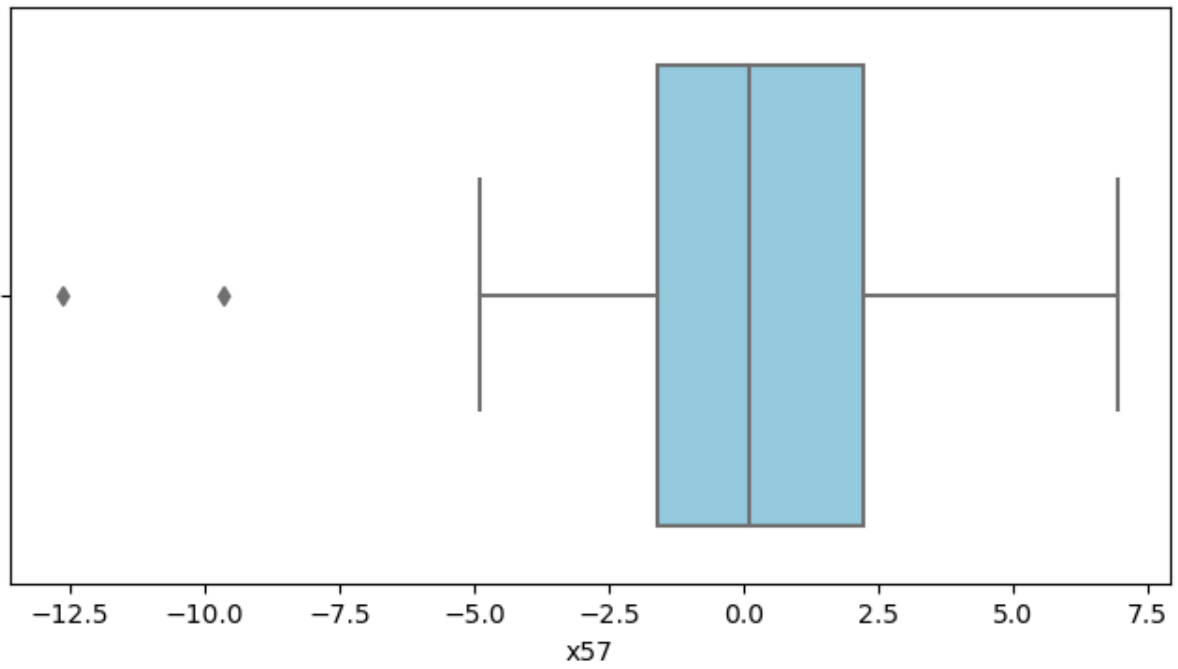
Boxplot for x55



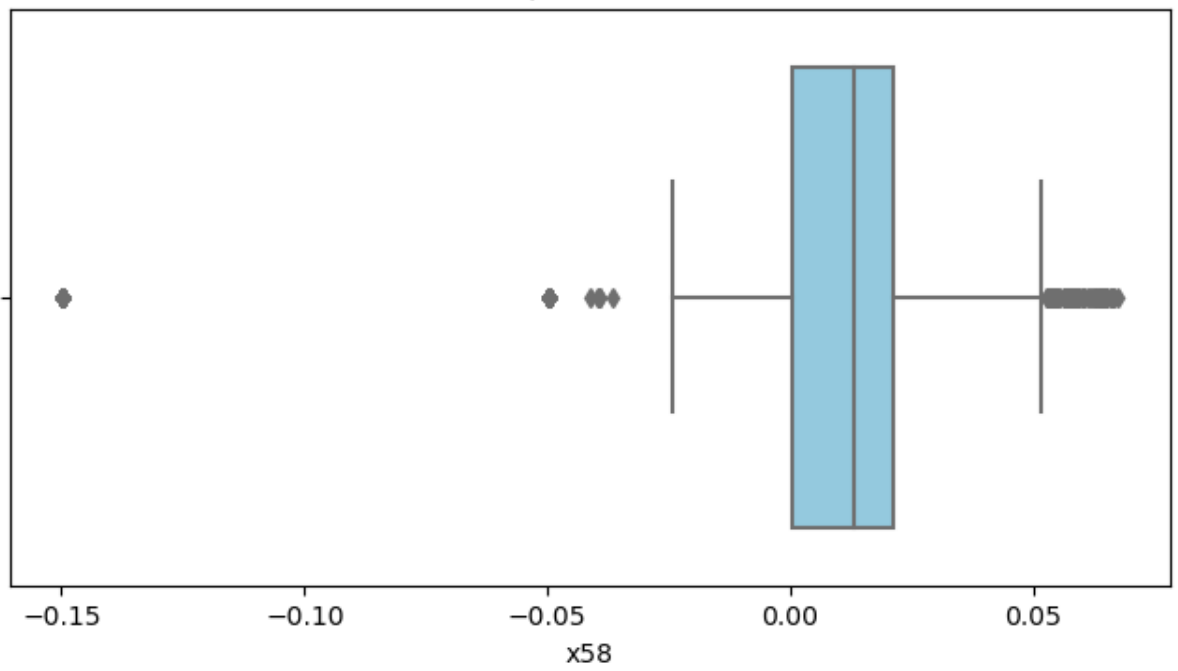
Boxplot for x56



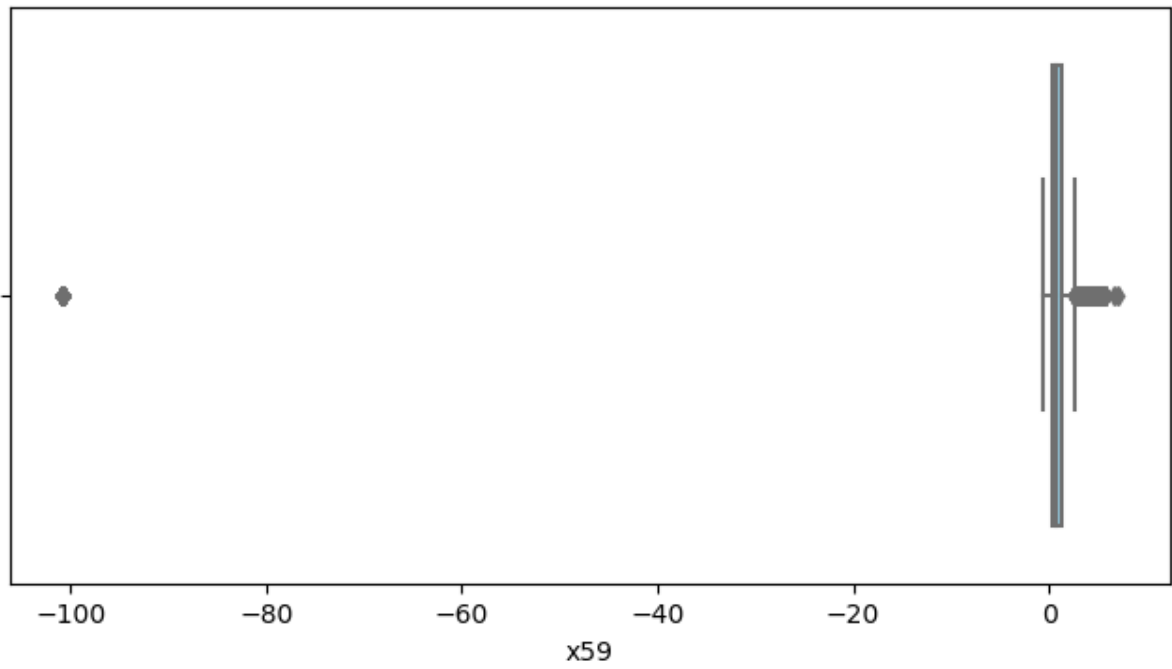
Boxplot for x57



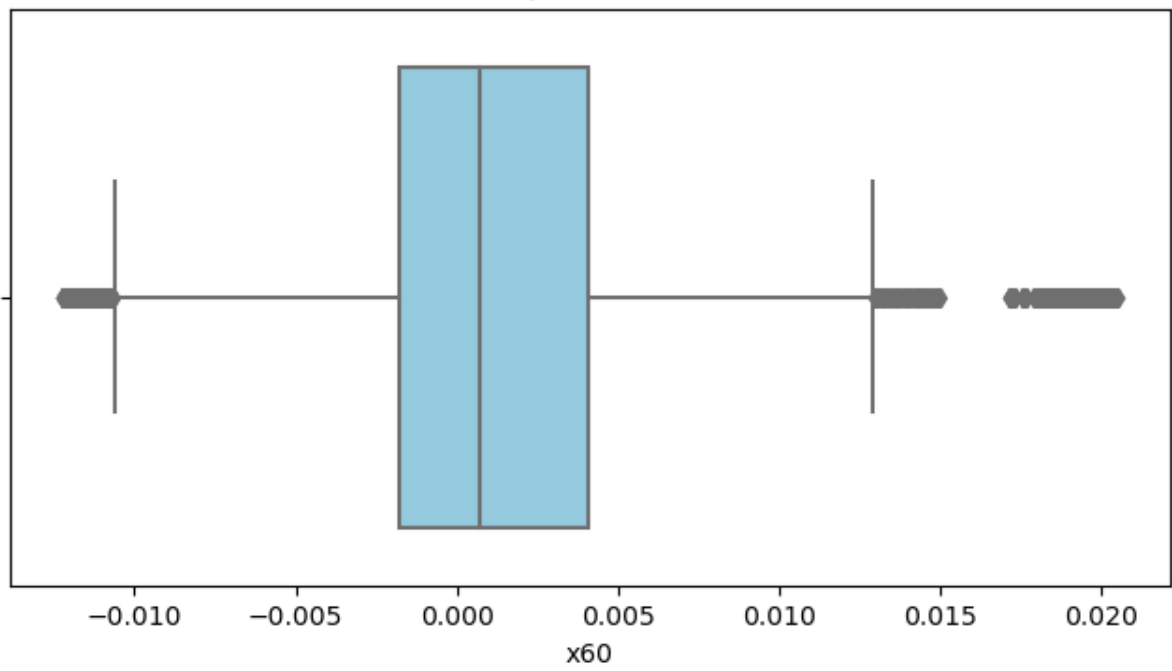
Boxplot for x58

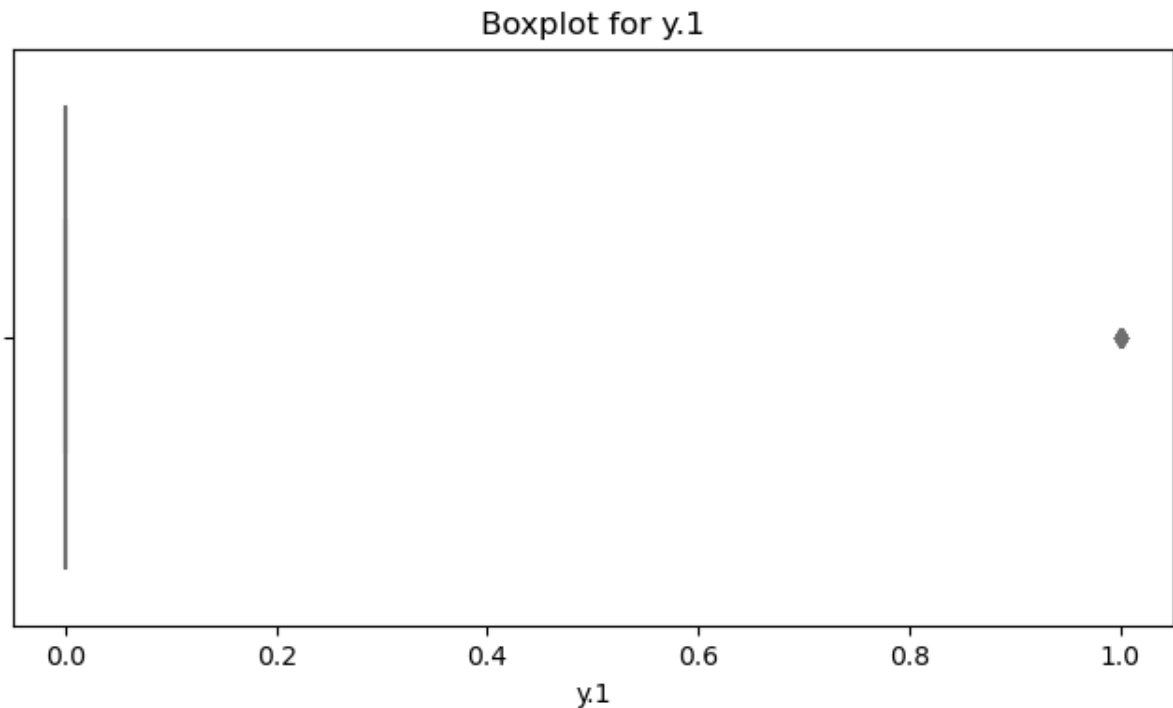


Boxplot for x59



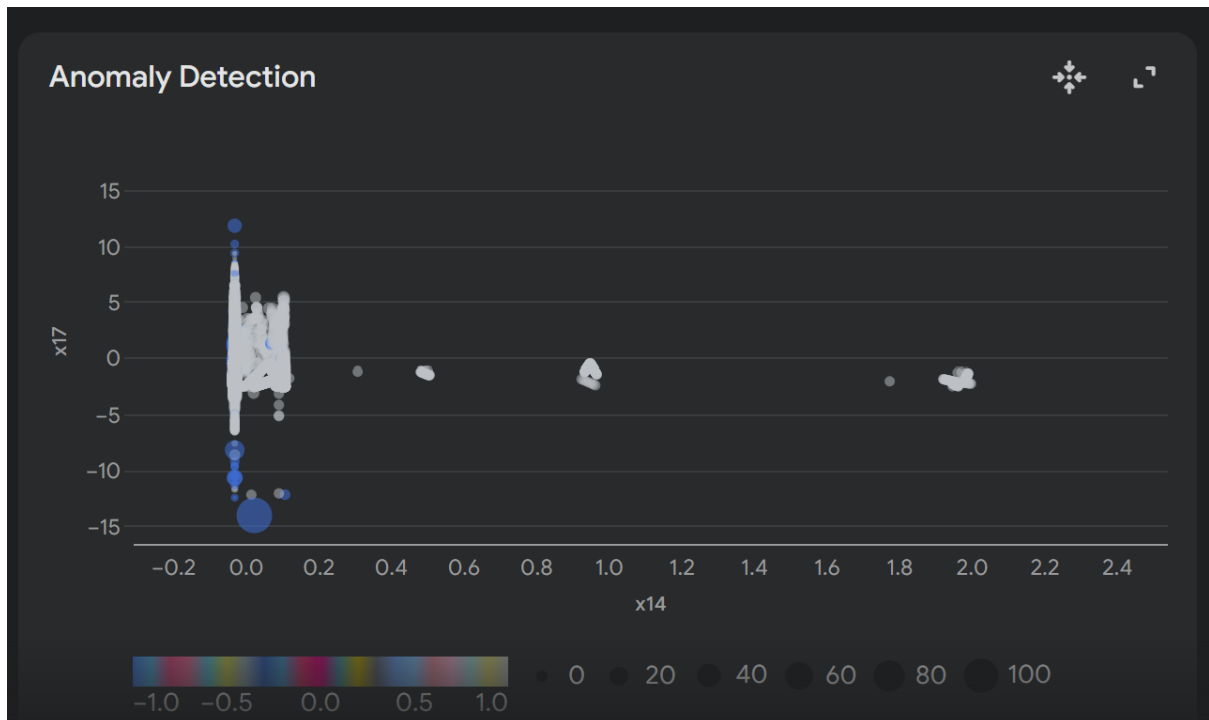
Boxplot for x60





3. Anomaly Detection with Isolation Forest

The Isolation Forest algorithm was chosen for anomaly detection due to its effectiveness in identifying anomalies in high-dimensional datasets. This algorithm isolates anomalies by randomly selecting a feature and then randomly creating a split value between the maximum and minimum values of the selected feature. The process is repeated recursively until all data points are isolated. Anomalies are identified as the points that are isolated first, requiring fewer splits.



4. Model Implementation and Evaluation

An Isolation Forest model was initialized with a `random_state` set to 42 for reproducibility. The model was trained on the entire dataset. Predictions were made on the same dataset to identify anomalies within it. The algorithm assigned a value of '1' to data points classified as normal and '-1' to those identified as anomalies.

5. Results and Discussion

The results of the anomaly detection process are stored in a new column named 'Predictions'. A '1' in this column indicates that the corresponding data point is not an anomaly, while a '-1' suggests it is an anomaly. The analysis successfully identified anomalies within the dataset.

6. Future Work

Future work may involve exploring other anomaly detection algorithms, such as One-Class SVM or Local Outlier Factor, and comparing their performance with the Isolation Forest. Additionally, fine-tuning the hyperparameters of the Isolation Forest, including the number of estimators and contamination rate, could potentially improve the accuracy of anomaly detection. Further investigation into the characteristics of the identified anomalies could provide valuable insights into the data and potential underlying issues.

7. Conclusion

The anomaly detection process using the Isolation Forest algorithm effectively identified unusual data points within the provided dataset. This information can be crucial for various data analysis tasks, such as data cleaning, fraud detection, and predictive modelling, by highlighting potential areas of concern or interest within the data.

This scatter plot visualizes the results of an Isolation Forest anomaly detection model. The model has identified two distinct groups within the dataset: 'normal' data points (represented by blue dots) and anomalies (represented by orange dots).

The majority of the data points have been classified as normal, forming a dense cluster towards the center of the visualization. The anomalies, on the other hand, appear scattered and isolated from the main cluster, indicating their deviation from the typical data patterns.

The size of the data points represents the magnitude of the values in the x24 column. Larger data points may indicate stronger deviations from normal patterns, further highlighting the outliers identified by the model.

Source Code

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.ensemble import IsolationForest

# Load the data
file_path = "AnomaData.xlsx" # Update with your file path if needed
data = pd.read_excel(file_path)
```

In [2]:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.ensemble import IsolationForest
from statsmodels.tsa.seasonal import seasonal_decompose
from statsmodels.tsa.stattools import adfuller
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
```

Cleaning The Dataset and Exploratory Data Analysis

```
# Display first 5 rows of the data
print("First 5 rows of the dataset:")
print(data.head())
```

```
# Basic info about dataset
print("\nDataset Information:")
print(data.info())
```

```

# Check for missing values
print("\nMissing values in each column:")
print(data.isnull().sum())

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.ensemble import IsolationForest
from statsmodels.tsa.seasonal import seasonal_decompose
from statsmodels.tsa.stattools import adfuller
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf

# Load the data
file_path = "AnomaData.xlsx" # Update with your file path if needed
data = pd.read_excel(file_path)

# Display first 5 rows of the data
print("First 5 rows of the dataset:")
print(data.head())

# Basic information about the dataset
print("\nDataset Information:")
data.info()

# Checking for missing values
print("\nMissing values in the dataset:")
print(data.isnull().sum())

# Drop columns with too many missing values (threshold = 50%)
missing_threshold = len(data) * 0.5
data = data.dropna(axis=1, thresh=missing_threshold)

# Fill remaining missing values with column mean for numerical columns
numerical_columns = data.select_dtypes(include=[np.number]).columns
data[numerical_columns] =
data[numerical_columns].fillna(data[numerical_columns].mean())

# Verify missing values are handled
print("\nMissing values after cleaning:")
print(data.isnull().sum())

# Descriptive statistics for numerical columns
print("\nDescriptive Statistics:")
print(data.describe())

# Check for duplicate rows
print("\nNumber of duplicate rows before cleaning:",
data.duplicated().sum())
data = data.drop_duplicates()
print("Number of duplicate rows after cleaning:", data.duplicated().sum())

# Distribution of target variable 'y'
if 'y' in data.columns:
    plt.figure(figsize=(8, 4))
    sns.countplot(x='y', data=data, palette='Set2')

```

```

plt.title("Distribution of Target Variable (y)")
plt.xlabel("y (1 = Anomaly, 0 = Normal)")
plt.ylabel("Count")
plt.show()

# Pairplot for understanding feature relationships
print("\nGenerating pairplot for numerical features:")
sns.pairplot(data, hue='y', palette='husl')
plt.show()

# Correlation matrix to understand relationships
print("\nCorrelation Matrix:")
plt.figure(figsize=(12, 8))
corr_matrix = data.corr()
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt='.2f')
plt.title("Correlation Matrix")
plt.show()

# Histograms for numerical columns
print("\nGenerating histograms for numerical features:")
numerical_columns = data.select_dtypes(include=[np.number]).columns
for col in numerical_columns:
    plt.figure(figsize=(8, 4))
    sns.histplot(data[col], kde=True, bins=30, color='blue')
    plt.title(f"Distribution of {col}")
    plt.xlabel(col)
    plt.show()

# Boxplots to identify outliers
print("\nGenerating boxplots to identify outliers:")
for col in numerical_columns:
    plt.figure(figsize=(8, 4))
    sns.boxplot(data[col], color='skyblue')
    plt.title(f"Boxplot of {col}")
    plt.show()

# Scatterplots for numerical features vs target
print("\nGenerating scatterplots for features against target variable 'y':")
if 'y' in data.columns:
    for col in numerical_columns:
        if col != 'y':
            plt.figure(figsize=(8, 4))
            sns.scatterplot(x=col, y='y', data=data, color='green')
            plt.title(f"Scatterplot of {col} vs y")
            plt.xlabel(col)
            plt.ylabel("y")
            plt.show()

# Time Series Check if 'timestamp' exists
print("\nChecking for time series analysis:")
if 'timestamp' in data.columns:
    data['timestamp'] = pd.to_datetime(data['timestamp'])
    data = data.set_index('timestamp')
    plt.figure(figsize=(12, 6))
    plt.plot(data.index, data['y'], color='tab:blue')
    plt.title("Time Series of Target Variable (y)")

```

```

plt.xlabel("Timestamp")
plt.ylabel("y")
plt.show()
else:
    print("No 'timestamp' column found for time series analysis.")

# Isolation Forest for anomaly detection
print("\nApplying Isolation Forest for Anomaly Detection:")
features = data.drop(columns=['y']) if 'y' in data.columns else data
iso_forest = IsolationForest(n_estimators=100, contamination=0.05,
random_state=42)
data['anomaly_score'] = iso_forest.fit_predict(features)
data['isolation_forest_anomaly'] = data['anomaly_score'].map({1: 0, -1: 1})

# Plot Isolation Forest results
plt.figure(figsize=(12, 6))
plt.plot(data.index, data['isolation_forest_anomaly'], color='red',
label="Anomalies")
plt.title("Isolation Forest Anomaly Detection")
plt.xlabel("Index / Timestamp")
plt.ylabel("Anomaly (1=Anomaly, 0=Normal)")
plt.legend()
plt.show()

# Summary of anomalies detected
print("\nSummary of Anomalies Detected by Isolation Forest:")
print(data['isolation_forest_anomaly'].value_counts())

Isolation Forest pipeline :
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import IsolationForest
from sklearn.metrics import classification_report, confusion_matrix
import joblib

# 1. Load and Preprocess the Data
def load_and_preprocess_data(file_path, target_column=None):
    data = pd.read_csv(file_path)
    data = data.dropna()

    if target_column:
        X = data.drop(columns=[target_column])
    else:
        X = data # Unsupervised case with no target

    scaler = StandardScaler()
    X_scaled = scaler.fit_transform(X)
    return X_scaled, scaler

# 2. Train Isolation Forest Model
def train_isolation_forest(X_train, contamination=0.1):
    model = IsolationForest(contamination=contamination, random_state=42)
    model.fit(X_train)
    return model

```

```

# 3. Evaluate the Model (if true labels exist)
def evaluate_model(model, X_val, y_true):
    y_pred = model.predict(X_val)
    y_pred = [1 if x == 1 else 0 for x in y_pred] # Convert -1 to 0
(outliers)
    print("Classification Report:")
    print(classification_report(y_true, y_pred))
    print("Confusion Matrix:")
    print(confusion_matrix(y_true, y_pred))

# 4. Save the Model and Scaler
def save_model(model, scaler, model_path='isolation_forest_model.pkl',
scaler_path='scaler.pkl'):
    joblib.dump(model, model_path)
    joblib.dump(scaler, scaler_path)

# 5. Load the Model and Make Predictions
def load_and_predict(model_path, scaler_path, new_data_path):
    model = joblib.load(model_path)
    scaler = joblib.load(scaler_path)

    new_data = pd.read_csv(new_data_path)
    new_data_scaled = scaler.transform(new_data)

    predictions = model.predict(new_data_scaled)
    predictions = [1 if x == 1 else 0 for x in predictions]
    return predictions

# 6. Exploratory Data Analysis for Anomalies
def exploratory_data_analysis(file_path):
    import matplotlib.pyplot as plt
    import seaborn as sns

    # Load the dataset
    data = pd.read_csv(file_path)

    print("Data Overview:")
    print(data.head())
    print("\nData Description:")
    print(data.describe())
    print("\nMissing Values:")
    print(data.isnull().sum())

    # Pair plot for understanding relationships
    print("\nGenerating pair plot...")
    sns.pairplot(data)
    plt.show()

    # Boxplots to detect outliers
    print("\nGenerating boxplots for each feature...")
    for column in data.columns:
        plt.figure(figsize=(10, 6))
        sns.boxplot(data[column])
        plt.title(f"Boxplot for {column}")
        plt.show()

    # Correlation heatmap

```

```

print("\nGenerating correlation heatmap...")
plt.figure(figsize=(12, 8))
sns.heatmap(data.corr(), annot=True, cmap='coolwarm')
plt.title("Feature Correlation Heatmap")
plt.show()

# Main Pipeline
if __name__ == "__main__":
    # File paths
    data_path = 'your_dataset.csv'
    new_data_path = 'new_data.csv'
    target_column = None # Set to your target column if labeled data
exists

    # Load and preprocess data
    X, scaler = load_and_preprocess_data(data_path, target_column)
    X_train, X_val = train_test_split(X, test_size=0.2, random_state=42)

    # Train model
    model = train_isolation_forest(X_train, contamination=0.1)

    # If labeled data is available for evaluation
    # y_true = ... # Provide true labels here
    # evaluate_model(model, X_val, y_true)

    # Save model and scaler
    save_model(model, scaler)

    # Predict on new data
    predictions = load_and_predict('isolation_forest_model.pkl',
'scaler.pkl', new_data_path)
    print("Predictions on new data:", predictions)

    # Perform Exploratory Data Analysis
    print("\nPerforming Exploratory Data Analysis...")
    exploratory_data_analysis(data_path)

```