## ▾ HW10 CS395 apfriedr 4/24/2023

What are the challenges associated with outlier detection in data analysis, and how do different types of outliers affect the process of identifying and managing outliers in a dataset? Provide technical reasoning to support your answer

Outlier detection is the process of identifying data points that are very different from expectation. There are several problems typical to outlier analysis or anomaly detection. outliers can be global or contextual - dependent on other data they come with, or completely and obviously apart in their value. There can also be collective outliers, such as orders through a day that are delayed. This delay may or may not be caused by the same mechanism, despite us identifyinng this group using the same methodology. Outlier detection depends on the dsitribution or modeling of non-outliers. Therefore if this modeling is done poorly or is very noise, it can very difficult to tell what is and what isn't an outlier, or to be able to trust the results of detection. Choosing the similarity/distance measure to detect outliers is often also application dependent, meaning some amount of domain knowledge is likely required to detect outliers. Finally, we may be able to detect outliers but also want to know \textit{why} they are outliers, which may not always be possible. We can use statistical methods to judge the probability of two outliers in relationship to each other, which may gives us some notion of whether they were generated by the same mechanism.

Supervised and unsupervised methods are both used to detect outliers. Supervised methods involve labeling and classifying whether data is likely to be part of the built model. Two challenges that can occur are data imbalance, and sometimes the recall is more important than the simple labelling. For unsupervised methods, clustering is often applied. This relies on an assumption that the data points are clustered inn some way to begin with. Data points not belonging to a cluster are marked as outliers. This can lead to fallacious results when using large datasets, unless small groups of indiidual points are identified as outlier clusters as well.

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import scipy.stats as stats
5 from scipy.stats import chi2
6 from numpy.linalg import inv
7 import seaborn as sns
```

```
1 np.random.seed(123)  # For reproducibility
2
3 # Generate synthetic multivariate dataset with outliers
```

● ✕

```python
5  n_features = 5  # Number of features
6  n_outliers = np.random.randint(100, 150)  # Number of outliers
7  mean = np.zeros(n_features)
8  cov = np.identity(n_features)
9  data = np.random.multivariate_normal(mean, cov, n_samples)
10 outliers = np.random.uniform(low=-10, high=10, size=(n_outliers, n_features))
11 data = np.concatenate((data, outliers), axis=0)
12 data = pd.DataFrame(data, columns=[f'Feature_{i+1}' for i in range(n_features)])
13 data.to_csv('synthetic_dataset.csv', index=False)
14
```

```python
1  # Load the synthetic dataset from the "synthetic_dataset.csv" file into a pandas DataFra
2  df = pd.read_csv("synthetic_dataset.csv", header=0)
3
4  print(df.head(10))
5
```

```
        Feature_1   Feature_2   Feature_3   Feature_4   Feature_5
0   -0.952097   -0.745441    1.738370    0.406454    0.322106
1   -0.051518   -0.204201    1.979348   -1.619300   -1.113964
2   -0.447441    1.668402   -0.143372   -0.619191   -0.769433
3    0.576746    0.126526   -1.301489    2.207427    0.522742
4    0.465645    0.724915    1.495827    0.746581   -1.100985
5   -1.410301   -0.747651   -0.984868   -0.748569    0.240367
6   -1.855637   -1.779455   -2.750224   -0.234158   -0.695981
7   -1.774134    2.361601    0.034993   -0.344642   -0.725032
8    1.039606   -0.241728   -0.112905   -1.660696    0.013539
9    0.337374   -0.926623    0.275747    0.370852    1.174307
```

```python
1  def grubbs_test(df, alpha):
2      # Calculate the mean and standard deviation for each column of the DataFrame
3      mean = df.mean()
4      std_dev = df.std()
5
6      # Calculate the test statistic for each column of the DataFrame
7      G = abs(df - mean) / std_dev
8      max_G = G.max(axis=1)
9
10     # Calculate the critical value from the t-distribution
11     n = df.shape[0]
12     t = stats.t.ppf(1 - alpha / (2 * n), n - 2)
13     critical_value = (n - 1) / np.sqrt(n) * np.sqrt(np.square(t) / (n - 2 + np.square(t)
14
15     # Return a Boolean mask indicating which rows are identified as outliers
16     return max_G > critical_value
17
```

```python
1  # Mahalanobis distance
```

```python
2 def calculateMahalanobis(y=None, data=None, cov=None):
3
4     y_mu = y - np.mean(data)
5     if not cov:
6         cov = np.cov(data.values.T)
7     inv_covmat = np.linalg.inv(cov)
8     left = np.dot(y_mu, inv_covmat)
9     mahal = np.dot(left, y_mu.T)
10     return mahal.diagonal()
11
12 def find_outliers_mahalanobis(df, mahal_dist_col, alpha=0.05):
13     """
14     Find outliers using Mahalanobis distance and a given significance level (alpha).
15
16     Parameters:
17     df (pandas.DataFrame): DataFrame with the Mahalanobis distances calculated in a colu
18     mahal_dist_col (str): The name of the column containing Mahalanobis distances.
19     alpha (float): Significance level to determine outliers (default is 0.05).
20
21     Returns:
22     pandas.Series: A boolean mask indicating whether each row is an outlier (True) or no
23     """
24
25     # Calculate degrees of freedom (number of variables/features)
26     df_features = df.drop(columns=[mahal_dist_col])
27     num_features = df_features.shape[1]
28
29     # Calculate Chi-squared critical value for the given alpha
30     chi2_critical_value = stats.chi2.ppf(1 - alpha, num_features)
31
32     # Create a boolean mask for outliers
33     outliers_mask = df[mahal_dist_col] > chi2_critical_value
34
35     return outliers_mask
```

```python
1 # Chi-square statistical approach.
2
```

```python
1 def plot_outliers_boxplot(df, outliers, test_name):
2     # Print the indices of the rows identified as outliers
3     print(df[outliers].index)
4
5     # Create a box plot of the non-outlier data
6     non_outliers = df[~outliers]
7     non_outliers.plot(kind='box')
8
9     # Plot the outliers as separate points
10     outlier_rows = df[outliers]
11     plt.plot(range(1, df.shape[1] + 1), outlier_rows.values.T, 'ro')
```
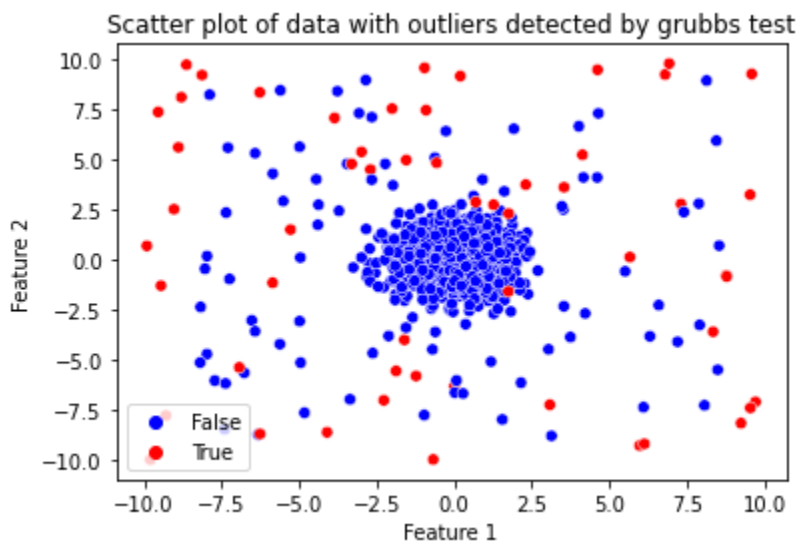
```
12
13      # Set the title and labels for the plot
14      plt.title('Box plot of data with outliers detected by {test_name} test'.format(test_
15      plt.xlabel('Features')
16      plt.ylabel('Values')
17
18      # Show the plot
19      plt.show()
```

```
 1 def plot_outliers_scatterplot(df, outliers, test_name):
 2      # Create a scatter plot of the data, coloring the outliers differently
 3      sns.scatterplot(data=df, x='Feature_1', y='Feature_2', hue=outliers, palette=['blue
 4
 5      # Set the title and labels for the plot
 6      plt.title('Scatter plot of data with outliers detected by {test_name} test'.format(t
 7      plt.xlabel('Feature 1')
 8      plt.ylabel('Feature 2')
 9
10      # Show the plot
11      plt.show()
```

```
 1 # Apply the Grubbs' test with a significance level of 0.05
 2 grubbs_outliers = grubbs_test(df, 0.05)
 3 plot_outliers_scatterplot(df, grubbs_outliers, "grubbs")
 4 plot_outliers_boxplot(df, grubbs_outliers, "grubbs")
```
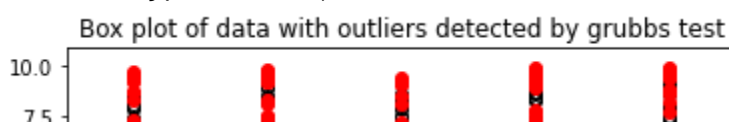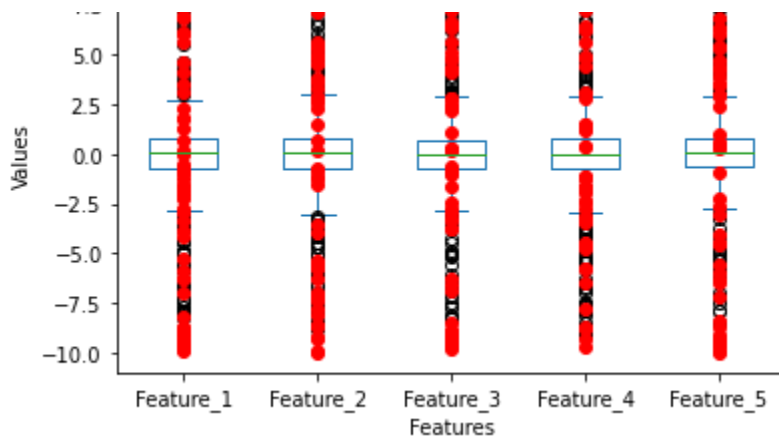


```
Int64Index([1000, 1002, 1009, 1012, 1019, 1020, 1023, 1026, 1028, 1032, 1034,
            1035, 1036, 1037, 1040, 1041, 1044, 1045, 1051, 1055, 1057, 1059,
            1060, 1068, 1069, 1070, 1071, 1073, 1074, 1075, 1076, 1078, 1081,
            1084, 1085, 1087, 1099, 1102, 1110, 1111, 1112, 1114, 1115, 1117,
            1120, 1124, 1125, 1127, 1128, 1129, 1131, 1133, 1134, 1135, 1140],
           dtype='int64')
```
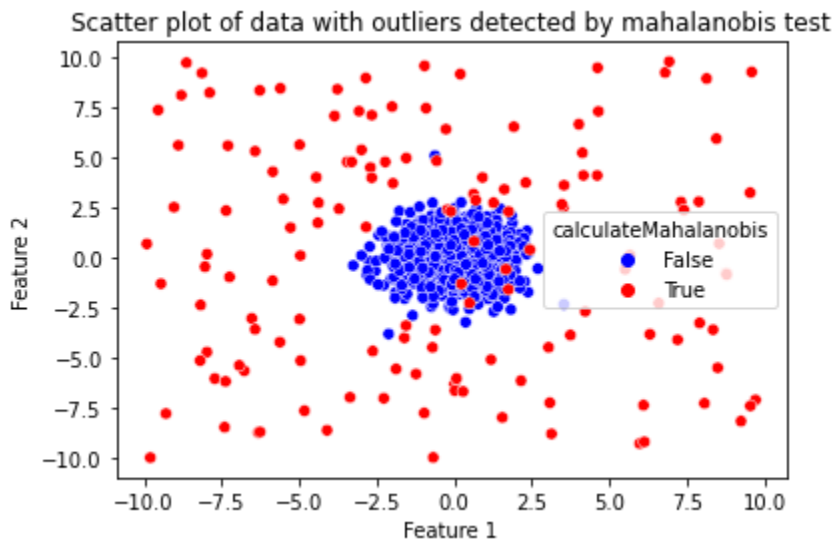
```
1 # Apply the Mahalanobis outlier detection with a significance level of 0.05
2 df_mahal = df.copy()
3 df_mahal['calculateMahalanobis'] = calculateMahalanobis(y=df, data=df[['Feature_1',  'Fe
4 # df['p'] = 1 - chi2.cdf(df['Mahalanobis'], 4)
5
6
7
```

```
1 mahal_outliers = find_outliers_mahalanobis(df_mahal, "calculateMahalanobis")
2 plot_outliers_scatterplot(df, mahal_outliers, "mahalanobis")
3 plot_outliers_boxplot(df, mahal_outliers, "mahalanobis")
```
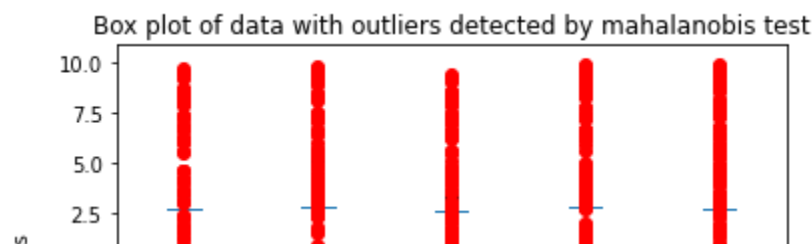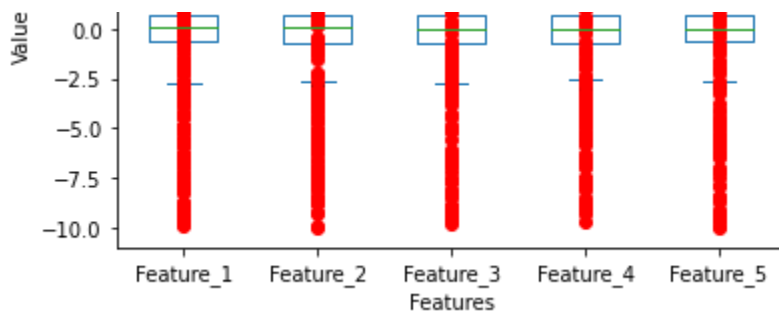


```
Int64Index([1000, 1001, 1002, 1003, 1004, 1005, 1006, 1007, 1008, 1009,
            ...
            1134, 1135, 1137, 1138, 1139, 1140, 1141, 1142, 1143, 1144],
           dtype='int64', length=141)
```

```
1 def find_outliers_chi_squared(df, alpha=0.05):
2     """
3     Find outliers using Chi-squared statistic approach on a DataFrame.
4
5     Parameters:
6     df (pandas.DataFrame): DataFrame containing numeric data.
7     alpha (float): Significance level to determine outliers (default is 0.05).
8
9     Returns:
10    pandas.Series: A boolean mask indicating whether each row is an outlier (True) or no
11    """
12    # Calculate z-scores for each value in the DataFrame
13    z_scores = np.abs(stats.zscore(df))
14
15    # Calculate the Chi-squared statistic for each row (sum of squares of z-scores)
16    chi2_stat = np.sum(z_scores**2, axis=1)
17
18    # Calculate degrees of freedom (number of variables/features)
19    num_features = df.shape[1]
20
21    # Calculate Chi-squared critical value for the given alpha
22    chi2_critical_value = stats.chi2.ppf(1 - alpha, num_features)
23
24    # Create a boolean mask for outliers
25    outliers_mask = chi2_stat > chi2_critical_value
26
27    return pd.Series(outliers_mask, index=df.index)
28
29
```
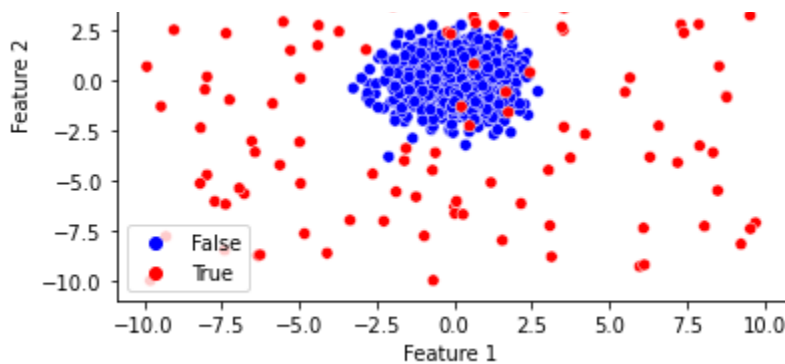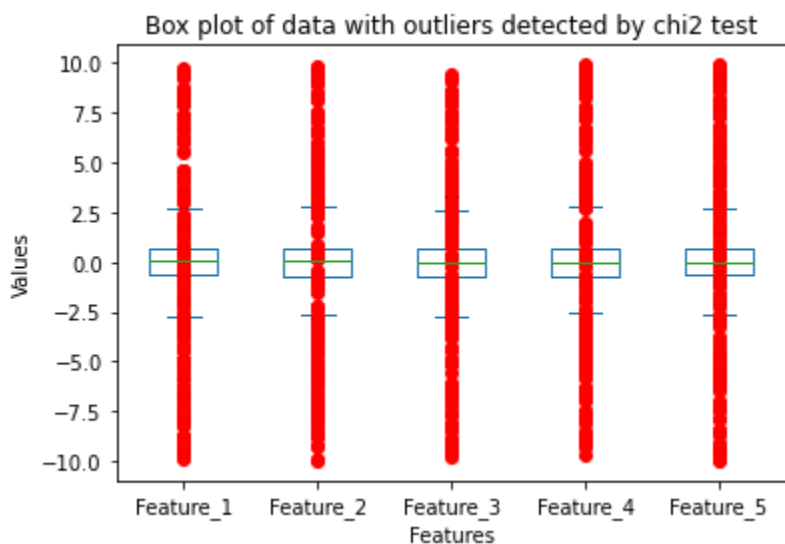
```
1 chi2_outliers = find_outliers_chi_squared(df)
2 plot_outliers_scatterplot(df, chi2_outliers, "chi2")
3 plot_outliers_boxplot(df, chi2_outliers, "chi2")
4
```

```
Int64Index([1000, 1001, 1002, 1003, 1004, 1005, 1006, 1007, 1008, 1009,
            ...
            1134, 1135, 1137, 1138, 1139, 1140, 1141, 1142, 1143, 1144],
           dtype='int64', length=142)
```



# Analysis of outlier detection methods

I plotted scatterplots of feature one vs feature 2 for each test detection method. It would be possible to plot the cartesian product of the features, but that amount of output seemed excessive - particularly given that visually, the detected outliers are mainly dependent on these first 2 features. The detection methods worked well, finding points mostly outside of the main feature 1 and 2 cluster. In the box plots, it can be seen that Grubbs' test did a better job of findinng outliers, particularly for features 3 and 5. If one wanted to go more in depth on this analysis, I would suggest to do a principal components analysis to aid in the visualization of clusters and outliers to those clusters. This would likely mitigate the problem of having many features, by reducing the dimensions of the data. Finally, comparing the actual outliers detected, it seems that all 3 methods detected many of the same outliers, while Grubbs' test only detected around 1/3 as many. We expect 100-150 outliers given our sythetic data, and both chi2 and Mahalanobis detected 142 and 141 respectively, while Grubbs' was about 55. This suggests that Grubbs' test was less reliable than either of the other 2 tests, or at least more stringent in what it

consdered an outlier.

Colab paid products  -  Cancel contracts here