

Outliers detection

Table of Contents

Boxplot.....	1
Local Distance-Based Outlier Factor (LDOF).....	9
Outliers deletion.....	11

On the basis of the majority voting method, we will apply first two methods for detecting outliers and then we will finally consider outliers just points classified by anomalous points by at least two of the algorithms implemented. These algorithms will be:

- Boxplot
- Local-Distance based outlier factor (LDOF)
- Local Outlier Factor

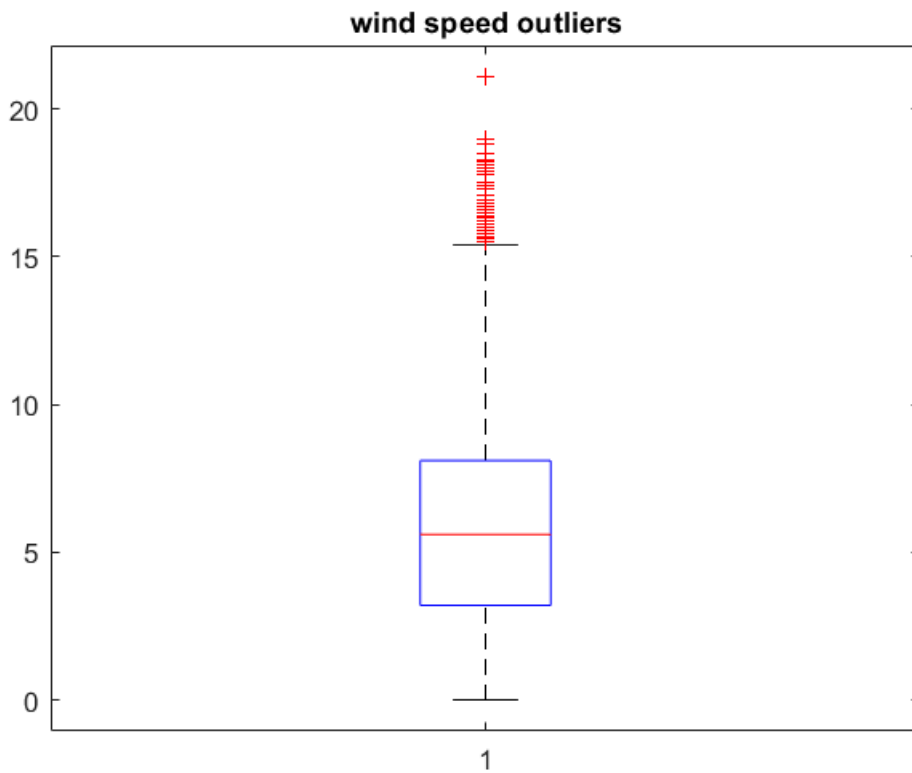
First, we put together all data

```
data10_Real = [data2010; data2011; data2012; data2013; data2014; data2015; data2016; data2017;
```

Boxplot

Before deleting outliers

```
clf;  
boxplot(data10_Real.WSPD)  
title('wind speed outliers')
```

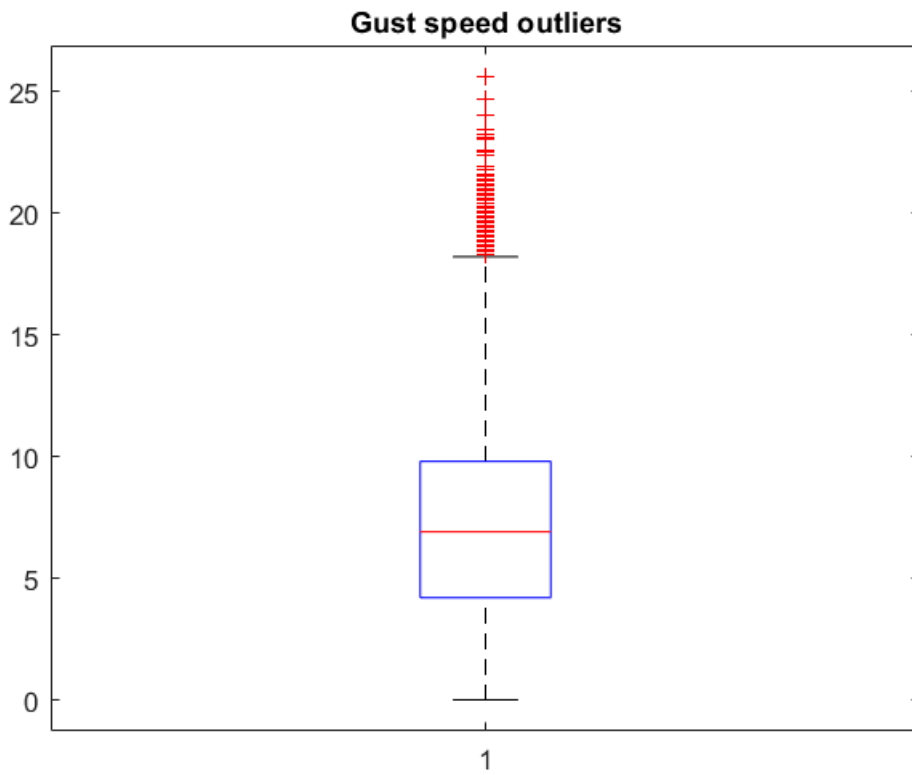


```
idx10OutWSPD = isoutlier(data10_Real.WSPD, 'quartiles')
```

```
idx10OutWSPD = 71031x1 logical array
```

```
0  
0  
0  
0  
0  
0  
0  
0  
0  
0  
0  
⋮
```

```
clf;  
boxplot(data10_Real.GST)  
title('Gust speed outliers')
```

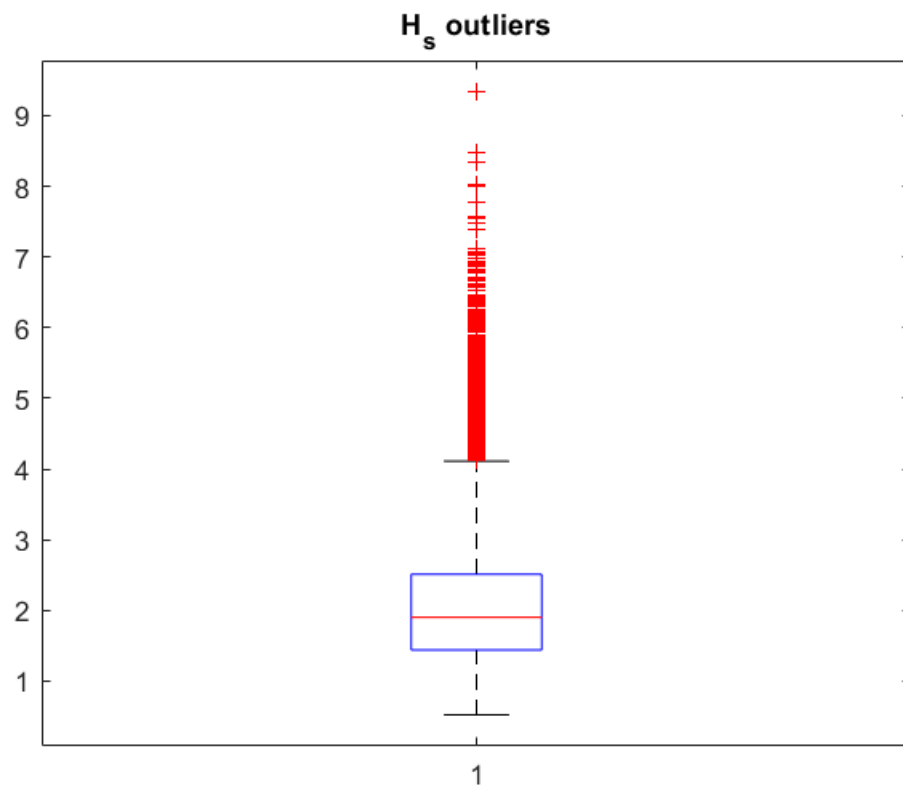


```
idx10OutGST = isoutlier(data10_Real.GST, 'quartiles')
```

```
idx10OutGST = 71031x1 logical array
```

```
0
0
0
0
0
0
0
0
0
0
0
⋮
```

```
clf;
boxplot(data10_Real.WVHT)
title('H_{s} outliers')
```

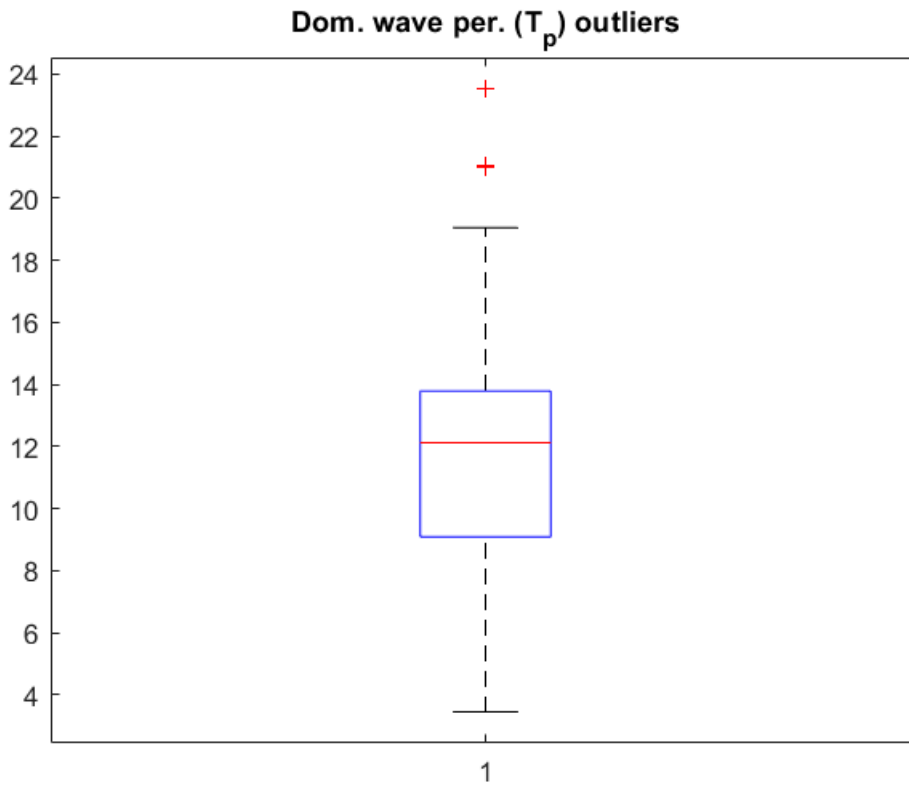


```
idx10OutWVHT = isoutlier(data10_Real.WVHT)
```

```
idx10OutWVHT = 71031x1 logical array
```

```
0
0
0
0
0
0
0
0
0
0
0
⋮
```

```
clf;
boxplot(data10_Real.DPD)
title("Dom. wave per. (T_{p}) outliers")
```

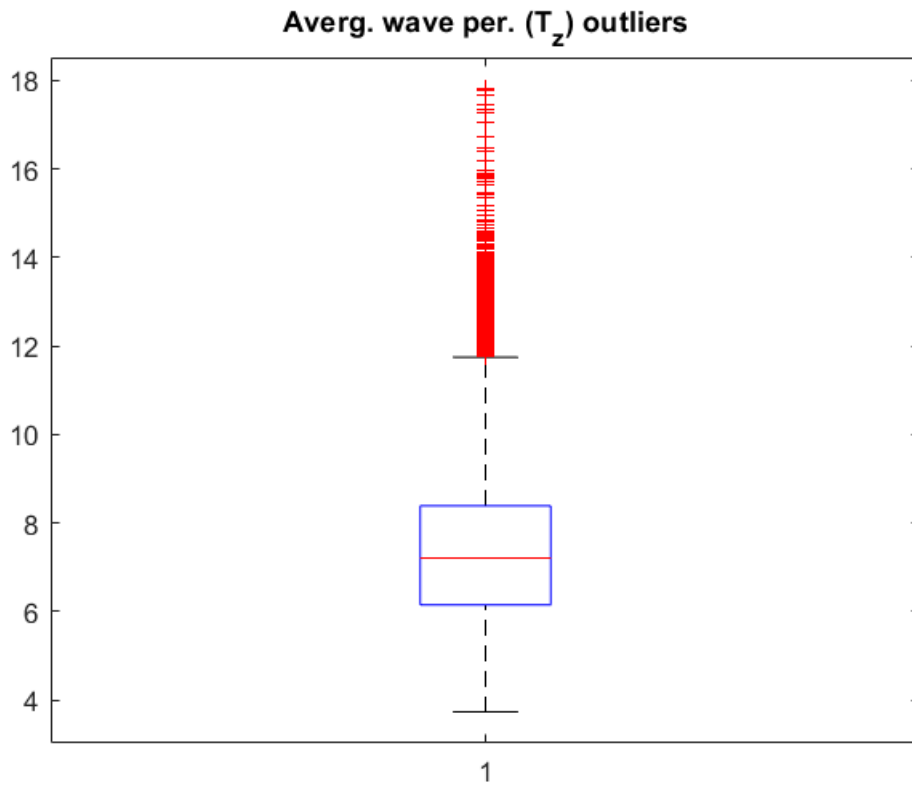


```
idx10OutDPD = find_outliers(data10_Real.DPD)
```

```
idx10OutDPD = 71031x1 logical array
```

```
0
0
0
0
0
0
0
0
0
0
0
⋮
```

```
clf;
boxplot(data10_Real.APD)
title('Averg. wave per. ( $T_z$ ) outliers')
```

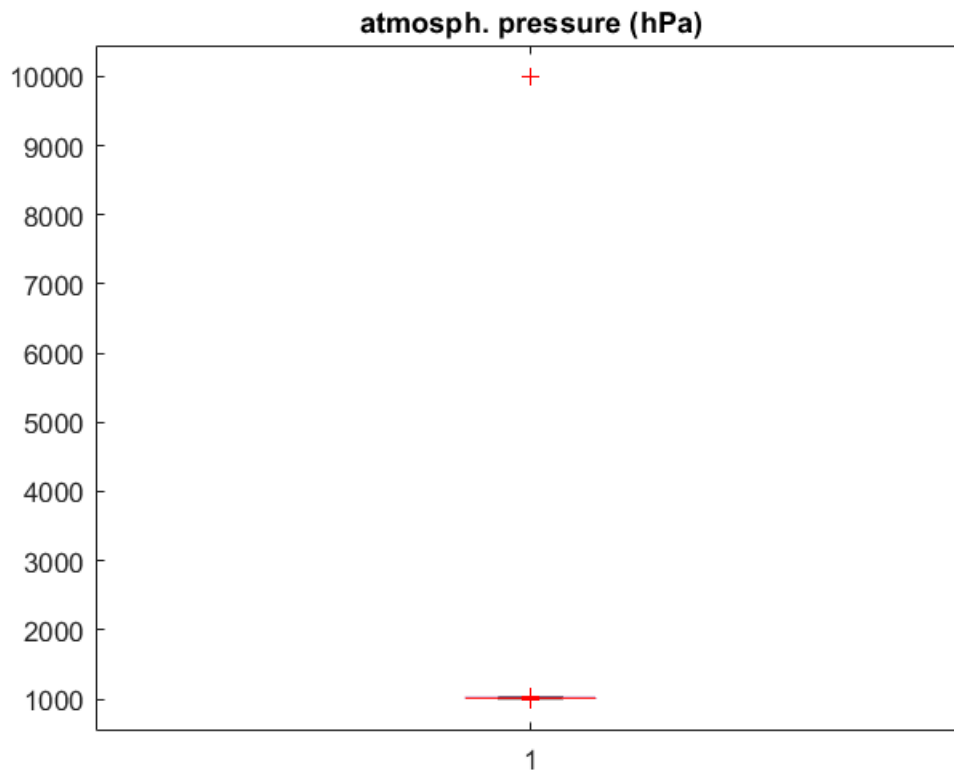


```
idx10OutAPD = find_outliers(data10_Real.APD)
```

```
idx10OutAPD = 71031x1 logical array
```

```
0
0
0
0
0
0
0
0
0
0
0
⋮
```

```
clf
boxplot(data10_Real.PRES)
title("atmosph. pressure (hPa)")
```

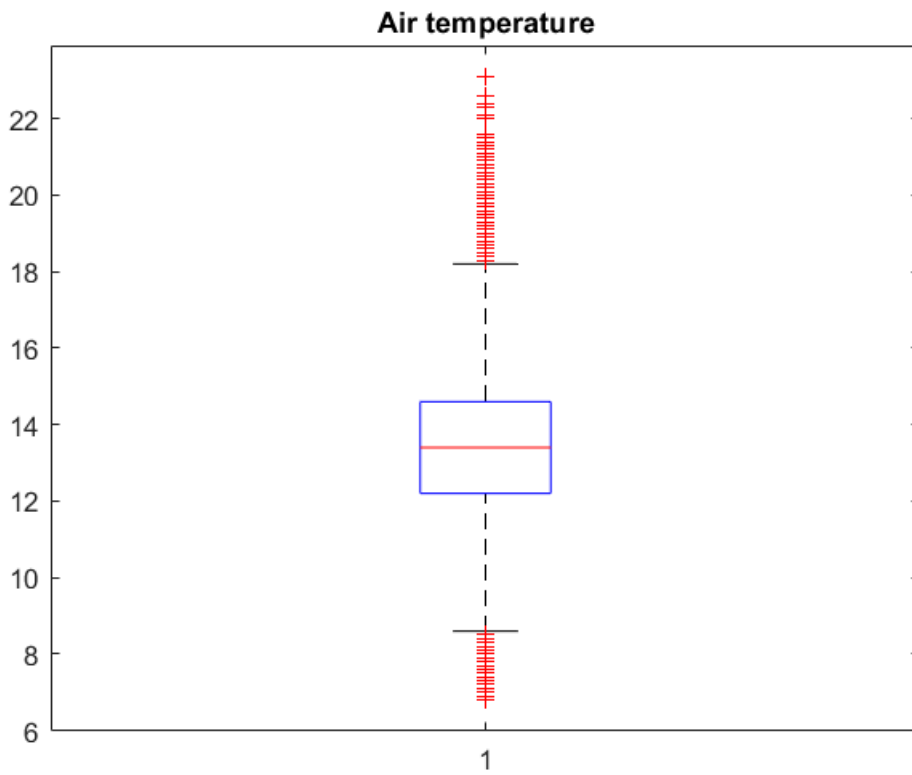


```
idx10OutPRES = find_outliers(data10_Real.PRES)
```

```
idx10OutPRES = 71031x1 logical array
```

```
0  
0  
0  
0  
0  
0  
0  
0  
0  
0  
0  
:  
:
```

```
clf  
boxplot(data10_Real.ATMP)  
title('Air temperature')
```

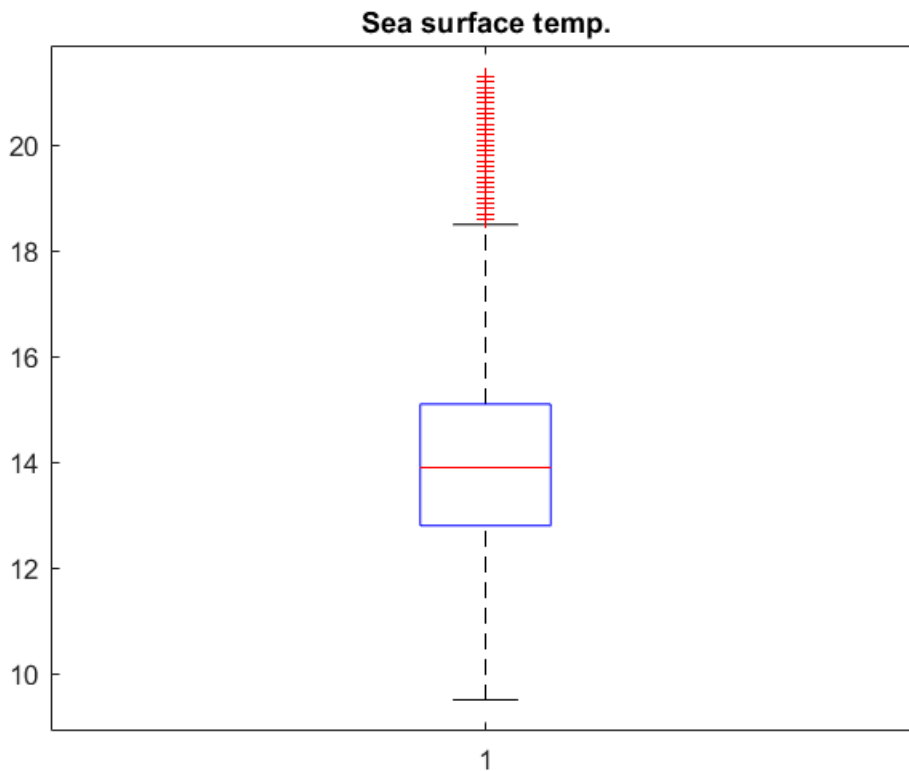


```
idx10OutATMP = find_outliers(data10_Real.ATMP)
```

```
idx10OutATMP = 71031x1 logical array
```

```
0  
0  
0  
0  
0  
0  
0  
0  
0  
0  
0  
⋮  
:
```

```
clf  
boxplot(data10_Real.WTMP)  
title("Sea surface temp.")
```

```
idx10OutWTMP = find_outliers(data10_Real.WTMP)
```

```
idx10OutWTMP = 71031x1 logical array
0
0
0
0
0
0
0
0
0
0
0
⋮
```

If we choose the complete data set for outliers detection, distribution measures from the station would be more realistic, because you aren't considering just one-year conditions, but you are taking points from several years and outliers would be anomaly points that don't occur almost any one of the years data.

Local Distance-Based Outlier Factor (LDOF)

Before applying this method

```
data10_Real = removevars(data10_Real, 'DATE');
```

```
[LDF0, D150,Idx150, KnnInnerDistance] = ldfo_outliers(table2array(data10_Real(:,{'WSPD','GST',
```

We have excluded Wind and Waves direction because they are circular variables. That means a value of 30° is the same as $30^\circ + 360^\circ$

```
LDF0
```

The ideal threshold would be 1, but it is dependent from problem so we have decided establish a threshold equals to the mean value of LDF0 coefficients of points.

```
threshold = mean(LDF0) %1.6069
```

```
threshold = 1.6069
```

```
idxOutlierLDF0 = (LDF0 > (threshold))
```

```
idxOutlierLDF0 = 71031x1 logical array
```

```
0
0
0
0
0
0
0
0
0
0
0
:
```

```
numFilasLDF0 = sum(idxOutlierLDF0)
```

```
idxOutlierBP = idx1OutWSPD | idx1OutWTMP | idx1OutATMP | idx1OutPRES | idx1OutAPD...
               | idx1OutDPD | idx1OutWVHT | idx1OutGST
```

```
idxOutlierBP = 71031x1 logical array
```

```
0
0
0
0
0
0
0
0
0
0
0
:
```

```
sum(idxOutlierBP)
```

```
ans = 5369
```

```
idxFinalOutlier = (idxOutlierBP == idxOutlierLDOF & idxOutlierBP == 1)
```

```
idxFinalOutlier = 71031x1 logical array  
0  
0  
0  
0  
0  
0  
0  
0  
0  
0  
0  
:  
:
```

```
sum(idxFinalOutlier)
```

We determine there are **1719** points with a high probability of being **outliers**, compared to 5330 detected by BoxPlot and 4820 detected by LDOF

Outliers deletion

```
data10_Real(idxFinalOutlier,:) = []
```