



## Chauvenet's Criterion



Chauvenet's criterion is a specific method for detecting outliers in data. It is based on the idea that, for a given dataset, the probability of an outlier occurring is relatively low. Chauvenet's criterion can be useful in certain situations, such as when you want to identify outliers in data that follows a normal distribution. However, it is not necessarily the best method for detecting outliers in all cases, and there are other methods that may be more appropriate in certain situations.

According to Chauvenet's criterion we reject a measurement (outlier) from a dataset of size  $N$  when its probability of observation is less than  $1/2N$ . A generalization is to replace the value 2 with a parameter  $C$ .

Source — Hoogendoorn, M., & Funk, B. (2018). Machine learning for the quantified self. *On the art of learning from sensory data*.

## Function

```
1 def mark_outliers_chauvenet(dataset, col, C=2):
2     """Finds outliers in the specified column of datatable and a
3     the same name extended with '_outlier' that expresses the re
4
5     Taken from: https://github.com/mhoogen/ML4QS/blob/master/Pyt
6
7     Args:
8         dataset (pd.DataFrame): The dataset
9         col (string): The column you want apply outlier detectio
10        C (int, optional): Degree of certainty for the identific
11                           of a normal distribution, typically be
12
13     Returns:
14         pd.DataFrame: The original dataframe with an extra boole
15                       indicating whether the value is an outlier or not.
16     """
17
18     dataset = dataset.copy()
19     # Compute the mean and standard deviation.
20     mean = dataset[col].mean()
21     std = dataset[col].std()
22     N = len(dataset.index)
23     criterion = 1.0 / (C * N)
24
25     # Consider the deviation for the data points.
26     deviation = abs(dataset[col] - mean) / std
27
28     # Express the upper and lower bounds.
29     low = -deviation / math.sqrt(C)
30     high = deviation / math.sqrt(C)
31     prob = []
32     mask = []
33
34     # Pass all rows in the dataset.
35     for i in range(0, len(dataset.index)):
36         # Determine the probability of observing the point
```

## Normal distribution

It's important to note that Chauvenet's criterion is only applicable to datasets that are normally distributed. If your dataset is not normally distributed, this method may not be suitable for identifying outliers.

- Histogram — Do you see a bell shaped curve?
- Boxplot — Is the box symmetrical?

