

SPCUP 2020: IEEE Signal Processing Cup 2020

Daniele Scapin, Gabriele Ferraresso, Marco Perin, Piero Simonetto, Riccardo Lorigiola

Abstract—In this letter, we proposed our method to solve the problem of unsupervised abnormality detection. It involves four main steps : peaks and abnormalities during data acquisition serch, anomaly check and system values update. The first and second steps take place at the same time during program flow but they serch autonomously the anomaly. Peaks serch is implemented with the use of polyfit and Kalman evaluation, for prediction of the measures and calculation of the percentual change. Instead abnormalities during data acquisition uses a modified version of the isolation forest to verify if during the measures are occurred discrepancies. The third part leverages of the concept that if every sensor in the autonomous system reports an anomaly, it means that the behaviour is not anomaly. The anomaly check uses a tree structure to get feedback from each sensor. The system values update is used to improve respose time and to update the variables needed in other sections.

Index Terms—Abnormalities detection, peaks search, Kalman filter, isoaltion forest, tree structure

I. INTRODUCTION

TO DO (Descrizione dettagliata flow del programma, associata al flow diagram)

II. PEAKS SEARCH: FINDPEAKSWRAPPER

FindPeaksWrapper(...) is used to find noticeable variation in the sensors data flow. The program can be divided in four main elements: polyfit, Kalman filter, peak presence and variables update.

A. Polyfit

For each sensor, for each new data acquisition the variation between the measure and the prediction is checked by *poly_fit*(...) function (it uses *polyfit*(...) and *polyval*(...) Matlab functions). This evaluation is done only if the number of elements to analyse is greater than degree of polyfit evaluation plus three (in the program degree plus three). This restriction is caused by the necessity of the presence of one element to verify and of degree plus two elements for *polyfit*(...)(plus two and not one because so the parametric estimation has al least a degree of freedom).

B. Kalman Filter

The Kalman filter is used only when the data are space, velocity and acceleration (angular or linear) and the acceleration is almost constant. The Kalman filter needs of the state transition matrix to be built (obtained from the differential

equations of the model). In this case the model is a black box model thus it can't be described exactly. For these reasons the Kalman filter is used only on space, velocity and acceleration with uniform acceleration. The Kalman filter is used anyway because if the autonomous system makes continuous and without variations paths/movements, it will generate a more accurate evaluation than *poly_fit*. The filter is activated by the evaluation of slope of regression line (the coefficient is generated inside *poly_fit* function). The covariance matrix ($Pn_{...}$) is initialized as $10 \text{ eye}(...)$ and update during subsequent cycles. The covariance matrix of the measures (Q) is set by *var2_error* vector placed diagonally (the motivation of these setting is explained in the next sections).

C. Peaks Presence

The peak presence is verified by calculating if the percentage of change between values expected and measured exceeds delta [%]. The value of delta is the maximum between a pre-set value (gap), the average percentage change of previous check and percentage error change of polyfit evaluation. This check is necessary because if signal noise is very loud, the function avoids reporting every measures as anomaly.

D. Internal variables update

Inside *FindPeaksWrapper*(...) there are three values (for each type of signal) that are used both as input and as output (Pn_2 , *varp_error*, *var2_error*) because in each cycle they are updated. Pn_2 is necessary for Kalman filter evaluation and in each cycle it is updated whitin the function itself. *varp_error* is average percentage change between measure and prediction (error) (related to prediction). *var2_error* is average squared change between measure and prediction (error) and it is used in the covariance matrix of the measure. It is used as covariance matrix of the measure because assuming the noises as gaussian noises (processes is subjected to ambient, thermal and internal noise)

$$\begin{aligned} y_{measure} &= y_{real} + e_{measure} & e_{measure} &\sim N(0, \Sigma_{measure}) \\ y_{estimate} &= y_{real} + e_{estimate} & e_{estimate} &\sim N(0, \Sigma_{estimate}) \end{aligned}$$

The measures of different axis are considered unrelated (covariance matrix is diagonal). Unidimensional estimation of variance

$$\begin{aligned} e &= y_{measure} - y_{estimate} & N(0, \theta) \\ e_i &= y_{measure_i} - y_{estimate_i} \\ l_e(\theta) &= -\log(\prod_{i=1}^n p_{e_i}(\theta)) \\ ... &= \frac{n}{2} \log(2\pi) + \frac{n}{2} \log(\theta) + \sum_{i=1}^n \frac{e_i^2}{2\theta} \\ \frac{\partial l_e(\theta)}{\partial \theta} &= \frac{n}{2\theta} - \sum_{i=1}^n \frac{e_i^2}{2\theta^2} \\ \hat{\theta} &= \frac{\sum_{i=1}^n e_i^2}{n} \end{aligned}$$

This competition is sponsored by the IEEE Signal Processing Society and MathWorks

All author is with the University of Padue, Padova, ITA

E. Input - Output

Detailed specification in *FindPeaksWrapper(...)*

Input

- t — time vector
- y — vector of values of $data_type$ element
- $data_type$ — type of data
- $degree$ — max degree during poly fit evaluation
- num — number of elements evaluating during polyfit
- ap — maximum permissible percentage error
- gap_sva — max variation to identify a constant element

Output

- $already_analysed$ — true if all values of the corresponding $data_type$ are already analysed
- $anomaly$ — vector of all the anomaly of the corresponding $data_type$
- $error$ — error of the last element ($y_{measured} - y_{predicted}$)
- y_next — prediction of the last element

III. ISOLATION FOREST

Isolation forest is a Statistical method to analyze small dataset and find anomalies. We used this algorithm to analyze data when the the points received are approximately constant. Isolation forest take a subsample of the data and ranomly select a dimension and a value in that dimension of the point, then points will be divided in two groups (points with a greater or equal value and points with lower value). By repeating this process until all points are isolated we are creating a binary tree. Repeating this process we obtain a forest of different trees where on the leaf of every tree we find the isolated points. Studing these trees we can obtain the average height of every point. Anomalies will have an average height lower than normal points. Now the anomaly score s of an instance x on a databes of n instances is defined as :

$$s(x, n) = 2^{\frac{E(h(x))}{c(n)}} \quad (1)$$

$$where \quad c(n) = 2H(n-1) - \frac{n-1}{n} \quad (2)$$

$$and \quad H(i) = \ln(i) + 0.5772156649 \quad (3)$$

A. Input

- NumTree: Number of trees inside the forest
- maxPoint: Maximum points inside every tree of the forest
- sk: Anomaly threshold
- type: Type of the new element
- newEl: New point of the forest

B. Output

- Last: True if newEl is abnormal, false otherwise
- Abnormal: True if there are abnormalities, False otherwise
- posOfAnomaly: Index of abnormal points
- h: Average height of each point in the forest
- s: Anomaly score of each point in the forest

IV. ANOMALY CHECK

This part contains classes and functions to structure data in a prioritized way, to avoid unnecessary controls over sensor that are less likely to give data representing a true anomaly

V. SYSTEM VALUES UPDATE

A. To do