



**SAPIENZA**  
UNIVERSITÀ DI ROMA

DEPARTMENT OF COMPUTER SCIENCE

# Highly Dimensional Anomaly Detection

## AUTOMATIC VERIFICATION OF INTELLIGENT SYSTEMS

**Professor:**

Enrico Tronci

**Students:**

Andrea Di Marco

Jemuel Espiritu Fuentes

{dimarco.1835169, fuentes.1803530}@studenti.uniroma1.it

---

Academic Year 2022/2023

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Implementation</b>	<b>3</b>
2.1	Geometric Entropy Minimization . . . . .	3
2.2	Principal Component Analysis . . . . .	5
2.3	Online Phase . . . . .	6
<b>3</b>	<b>Hyper-Parameter Optimization</b>	<b>7</b>
3.1	Geometric Entropy Minimization . . . . .	7
3.2	Principal Component Analysis . . . . .	8
<b>4</b>	<b>Testing</b>	<b>11</b>
4.1	Increasing Dimension . . . . .	11
4.2	Increasing Offline Samples . . . . .	12
4.3	Increasing Variance and Perturbation . . . . .	12
<b>5</b>	<b>Conclusions</b>	<b>14</b>
<b>6</b>	<b>Future Developments</b>	<b>14</b>
	<b>References</b>	<b>15</b>

# 1 Introduction

Anomaly detection is a key element of many real world applications like power networks or financial transactions and it is a vital aspect in both the security and safety of critical systems [1].

Many applications are relying on model based algorithms but in case of highly dimensional data streams, especially in real time systems, this is not a practical solution. An alternative is based on the extraction of useful univariate summary statistics from the observed data in order to detect anomalies.

For our course project we implemented an online anomaly detection system as proposed by [2] using `C++` and `Python` which works by computing summary statistics in an offline setting to be used as a comparison with real time data. The high dimensionality is mitigated by the use of two different algorithms depending on the use case. We subsequently performed optimization on the system in order to infer the most suitable values for the two main hyper-parameters  $\alpha$  and  $h$  on a number of different situations. Ulterior tests comparing the two models were also a part of our activities in order to observe the sensitivity of various statistics such mean and the dimensionality of the sample space. Lastly we discuss additional improvements and show some key weaknesses in the algorithms.

## 2 Implementation

We decided to implement the proposed anomaly detection algorithm in `C++` and interface it with `Python` during the optimization and testing phase and we used the `Eigen` library to facilitate linear algebra operations.

The system which is divided in an offline model training component and an online detection phase is shown in 1. In the offline phase the summary statistics are computed based on nominal data. In the online phase the summary statistics are compared to new data point to detect if an anomaly is present. The online phase utilizes a CUSUM-like algorithm to decide if the amount of difference between baseline and new data point is suspicious. The threshold for the alarm is represented by a hyperparameter  $h$ .

The computation of summary statistics, in both offline and online phase, is done by using a GEM or a PCA based approach depending on the use case, specifically the intrinsic dimensionality of the datasets.

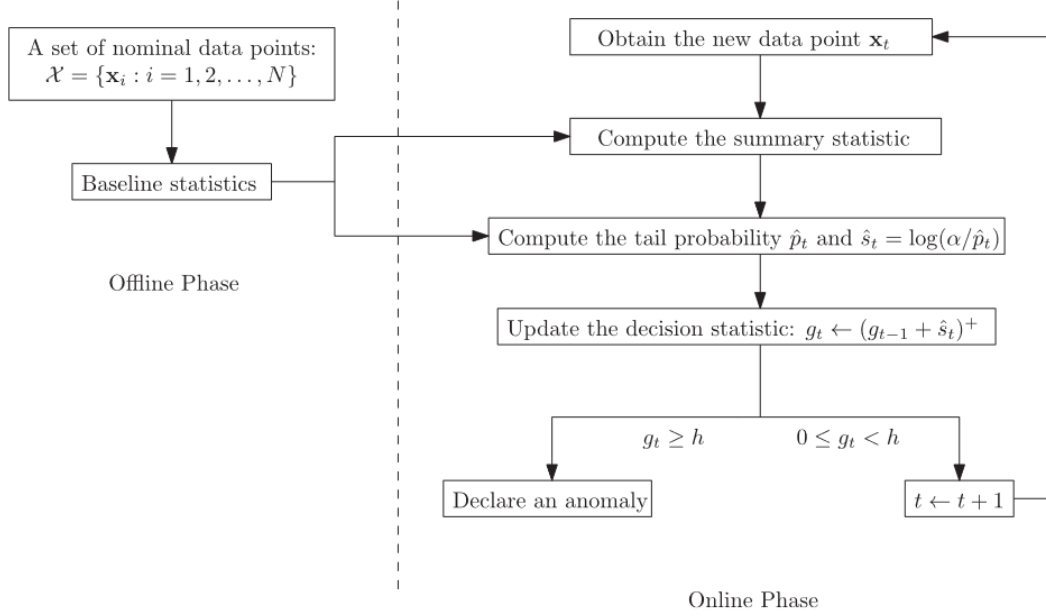


Figure 1: System diagram from [2]

### 2.1 Geometric Entropy Minimization

The first summary statistic is based on Geometric Entropy Minimization in which an acceptance region  $A$  is determined such that only outliers sit outside of  $A$  with respect to  $\alpha$ . The figure 2 shows the pseudocode of how the GEM based summary statistics is obtained. In this method first the dataset  $X$  is uniformly randomly partitioned in  $S_1$  and  $S_2$  with sizes  $N_1$  and  $N_2 = |X| - N_1$ . Then  $\forall x_j \in S_2$  the kNNs of  $x_j$  among each element of the set  $S_1$  is computed.

Let  $e_j(i)$  be the euclidean distance of  $x_j$  to its  $i$ th NN in  $S_1$ , we can represent the

sum of the distances from  $x_j$  to its  $k$  nearest neighbours as:

$$d_j \triangleq \sum_{i=1}^k e_j(i)$$

After computing  $\{d_j : x_j \in S_2\}$ , this set is sorted in ascending order and the acceptance region A is the  $(1 - \alpha)$  fraction of  $x_j \in S_2$  corresponding to the smallest  $(1 - \alpha)$  fraction of  $d_j$ 's. Then a new data point  $\mathbf{x}_t$  is considered an outlier if its sum of distances to its  $k$ NNs among  $S_1$ , denoted with  $d_t$ , is greater than the smallest  $(1 - \alpha)$  fraction of  $d_j$ 's which can be expressed as

$$\frac{\sum_{x_j \in S_2} \mathbf{1}\{d_t > d_j\}}{N_2} > (1 - \alpha)$$

where  $\mathbf{1}$  is the characteristic function.

In short, if  $\mathbf{x}_t$  is an outlier it falls out of the acceptance region A which means that the NN statistics  $d_t$  takes an higher value compared to non outliers. Moreover if the observed data stream persistently falls outside the acceptance region this could be an anomaly.

---

**Algorithm 1.** GEM-Based Real-Time Nonparametric Anomaly Detection

---

**Offline Phase**

- 1: Uniformly randomly partition the nominal dataset  $\mathcal{X}$  into two subsets  $S_1$  and  $S_2$  with sizes  $N_1$  and  $N_2$ , respectively.
- 2: **for**  $j : \mathbf{x}_j \in S_2$  **do**
- 3:   Search for the  $k$ NNs of  $\mathbf{x}_j$  among the set  $S_1$ .
- 4:   Compute  $d_j$  using (6).
- 5: **end for**
- 6: Sort  $\{d_j : \mathbf{x}_j \in S_2\}$  in ascending order.

**Online Detection Phase**

- 1: Initialization:  $t \leftarrow 0, g_0 \leftarrow 0$ .
  - 2: **while**  $g_t < h$  **do**
  - 3:    $t \leftarrow t + 1$ .
  - 4:   Obtain the new data point  $\mathbf{x}_t$ .
  - 5:   Search for the  $k$ NNs of  $\mathbf{x}_t$  among the set  $S_1$  and compute  $d_t$  using (6).
  - 6:    $\hat{p}_t = \frac{1}{N_2} \sum_{\mathbf{x}_j \in S_2} \mathbb{1}\{d_j > d_t\}$ .
  - 7:    $\hat{s}_t = \log(\alpha / \hat{p}_t)$ .
  - 8:    $g_t \leftarrow \max\{0, g_{t-1} + \hat{s}_t\}$ .
  - 9: **end while**
  - 10: Declare an anomaly and stop the procedure.
- 

Figure 2: Pseudocode of GEM based method from [2]

## 2.2 Principal Component Analysis

The second summary statistic is computed by using Principal Component Analysis. This particular algorithm was proposed by considering the curse of dimensionality which frequently affects real case datasets. In fact high dimensional dataset frequently exhibit a sparse structure and a lower dimension representation is feasible. In this case the model is effectively:

$$\mathbf{x}_t = \mathbf{y}_t + \mathbf{r}_t$$

where  $\mathbf{y}_t$  is the submanifold representation and  $\mathbf{r}_t$  is the residual term, mostly noise. The magnitude of the residual term can be used to characterize anomalous data because  $\|\mathbf{r}_t\|$  is higher for anomalies compared to nominal data. The figure 3 shows the steps in order to compute the summary statistics using a PCA based approach.

Given  $\mathbf{X}$  we can take two subsets  $S_1$  and  $S_2$  with sizes  $N_1$  and  $N_2$  both less than or equal to  $|\mathbf{X}|$ . Using  $S_1$  we can determine a submanifold where the data is embedded and using  $S_2$  we can compute the residual term  $R = \{\|\mathbf{r}_j\| : \mathbf{x}_j \in S_2\}$ .

First the sample mean  $\bar{\mathbf{x}}$  of  $S_1$  and the sample data covariance matrix  $\mathbf{Q}$  are computed.

Then the eigenvalues  $\{\lambda_j : j = 1, 2, \dots, p\}$  and the eigenvectors  $\{\mathbf{v}_j : j = 1, 2, \dots, p\}$  of  $\mathbf{Q}$  are obtained.

We can determine the dimensionality of the submanifold  $r$  based on the desired fraction of data variance retained given by

$$\gamma \triangleq \frac{\sum_{j=1}^r \lambda_j}{\sum_{j=1}^p \lambda_j} \leq 1$$

where the  $r$ -dimensional principal subspace is spanned by the orthonormal eigenvectors  $\mathbf{V} \triangleq \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r\}$  corresponding to the  $r$  largest eigenvalues of  $\mathbf{Q}$ .

The representation of  $\mathbf{x}_t$  in the linear submanifold can be determined as

$$\mathbf{y}_t = \bar{\mathbf{x}} + \sum_{j=1}^r \mathbf{v}_j \mathbf{v}_j^T (\mathbf{x}_t - \bar{\mathbf{x}}) = \bar{\mathbf{x}} + \mathbf{V} \mathbf{V}^T (\mathbf{x}_t - \bar{\mathbf{x}})$$

Then the residual term is easily computed by

$$\mathbf{r}_t = \mathbf{x}_t - \mathbf{y}_t = (\mathbf{I}_p - \mathbf{V} \mathbf{V}^T) (\mathbf{x}_t - \bar{\mathbf{x}})$$

. The summary statistic is  $\|\mathbf{r}_t\|_2$ .

---

**Algorithm 2.** PCA-Based Real-Time Nonparametric Anomaly Detection

---

**Offline Phase**

- 1: Choose subsets  $\mathcal{S}_1$  and  $\mathcal{S}_2$  of  $\mathcal{X}$  with sizes  $N_1$  and  $N_2$ , respectively.
- 2: Compute  $\bar{\mathbf{x}}$  and  $\mathbf{Q}$  over  $\mathcal{S}_1$  using (8) and (9), respectively.
- 3: Compute the eigenvalues  $\{\lambda_j : j = 1, 2, \dots, p\}$  and the eigenvectors  $\{\mathbf{v}_j : j = 1, 2, \dots, p\}$  of  $\mathbf{Q}$ .
- 4: Based on a desired level of  $\gamma$  (see (10)), determine  $r$  and form the matrix  $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r]$ .
- 5: **for**  $j : \mathbf{x}_j \in \mathcal{S}_2$  **do**
- 6:    $\mathbf{r}_j = (\mathbf{I}_p - \mathbf{V}\mathbf{V}^T)(\mathbf{x}_j - \bar{\mathbf{x}})$ .
- 7:   Compute  $\|\mathbf{r}_j\|_2$ .
- 8: **end for**
- 9: Sort  $\{\|\mathbf{r}_j\|_2 : \mathbf{x}_j \in \mathcal{S}_2\}$  in ascending order.

**Online Detection Phase**

- 1: Initialization:  $t \leftarrow 0, g_0 \leftarrow 0$ .
  - 2: **while**  $g_t < h$  **do**
  - 3:    $t \leftarrow t + 1$ .
  - 4:   Obtain the new data point  $\mathbf{x}_t$ .
  - 5:    $\mathbf{r}_t = (\mathbf{I}_p - \mathbf{V}\mathbf{V}^T)(\mathbf{x}_t - \bar{\mathbf{x}})$  and compute  $\|\mathbf{r}_t\|_2$ .
  - 6:    $\hat{p}_t = \frac{1}{N_2} \sum_{\mathbf{x}_j \in \mathcal{S}_2} \mathbb{1}\{\|\mathbf{r}_j\|_2 > \|\mathbf{r}_t\|_2\}$ .
  - 7:    $\hat{s}_t = \log(\alpha / \hat{p}_t)$ .
  - 8:    $g_t \leftarrow \max\{0, g_{t-1} + \hat{s}_t\}$ .
  - 9: **end while**
  - 10: Declare an anomaly and stop the procedure.
- 

Figure 3: Pseudocode of PCA based method from [2]

### 2.3 Online Phase

The summary statistics that are considered are  $\mathbf{d}_t$  for the GEM based technique and  $\|\mathbf{r}_t\|_2$  for the PCA technique. Outliers have statistics that are higher value compared to nominal data points. This means that outliers correspond to the right tail events considering the nominal pdfs of the summary statistics.

If  $\hat{p}_t < \alpha$  then we can claim that  $\mathbf{x}_t$  is an outlier with respect to the level of  $\alpha$ . Let:

$$s_t \triangleq \log\left(\frac{\alpha}{\hat{p}_t}\right)$$

Then we have  $s_t > 0$  if  $\mathbf{x}_t$  is an outlier and  $s_t \leq 0$  if it is a non outlier. The model free algorithm based on CUSUM therefore is:

$$\Gamma = \inf\{t : g_t \geq h\}$$

$$g_t = \max\{0, g_{t-1} + \hat{s}_t\}$$

$$g_0 = 0$$

### 3 Hyper-Parameter Optimization

We decided to optimize the two hyperparameters *alpha* and *h* of our models with *nevergrad*[3]. We utilized [3] *NGOpt* as the Black Box Optimization algorithm with the following objective function (1) where *a* is for the number of samples identified as anomalies in the training set,  $N_1$  is the amount of samples in the training set, *n* is the amount of samples identified as normal in the anomalous set,  $N_2$  the amount of samples in the anomalous set,  $w_1 = 0.4$ ,  $w_2 = 0.6$ .

$$l = w_1 \frac{a}{N_1} + w_2 \frac{n}{N_2} \quad (1)$$

We chose this specific function as a way to limit the possibility of the optimizer to simply pick very extreme thresholds, such as the case with only evaluating the fraction of anomaly detected by the system. The tests we run are the following.

1. **Change in the mean:** We trained the models on a 10 dimensional  $\mathcal{N} \sim (0, 1)$  distribution of 1.000 samples and tested it against a 10 dimensional  $\mathcal{N} \sim (1, 1)$  distribution of 10.000 samples.
2. **Change in the variance:** We trained the models on a 10 dimensional  $\mathcal{N} \sim (0, 1)$  distribution of 1.000 samples and tested it against a 10 dimensional  $\mathcal{N} \sim (0, 2)$  distribution of 10.000 samples.
3. **Subset:** We trained the models on a 10 dimensional  $\mathcal{N} \sim (0, 1)$  distribution of 1.000 samples and tested it against a 10 dimensional  $\mathcal{N} \sim (0.5, 0.5)$  distribution of 10.000 samples.
4. **Human movement:** We trained the models on a dataset representing embeddings of human activities [4] and tested it against anomalous data from the same dataset.

#### 3.1 Geometric Entropy Minimization

In the table 1 there can be seen the proposed hyperparameters and the resulting best achieved loss for the three test cases. In the images 5, 6 and 7 it can be observed the graph with the values assumed through time by  $\alpha$ , *h* and the *loss*. It is interesting to see how in 6 with the mean shift test case the algorithm reaches a plateau that it is unable to pierce.

It is worth noting how the GEM algorithm generally took significantly longer than the PCA algorithm, we believe it is due to GEM's computation of the *k* best neighbors amongst a partition of the training dataset. Said partition is currently set to be 15% of the original dataset, by reducing it the algorithm's speed will increase at the cost of its accuracy.



Test Case	$\alpha$	$h$	Loss	Time
$\mathcal{N} \sim (0.5, 0.5)$	0.036	5.290	0.6	1937s
$\mathcal{N} \sim (1, 1)$	0.019	8.124	0.6	1134s
$\mathcal{N} \sim (0, 2)$	0.535	1.001	0.206	1200s
<i>Human Activity</i>	0.0407	1.150	0.007	63931s

Table 1: Optimization results for GEM

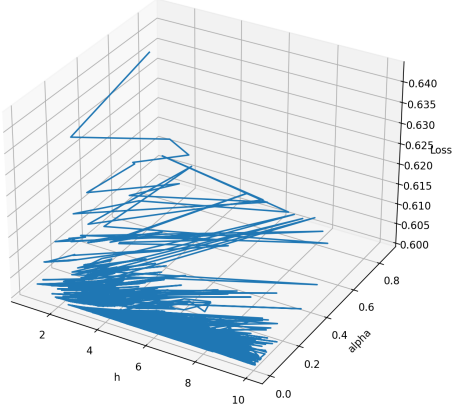


Figure 4: BBO of the GEM model with the subset test

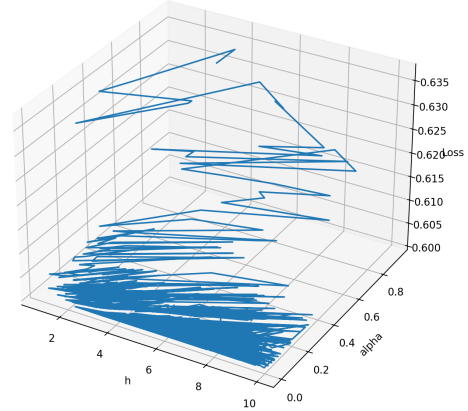


Figure 6: BBO of the GEM model with shifted mean test

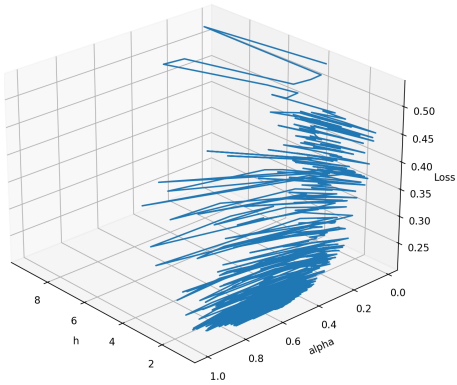


Figure 5: BBO of the GEM model with changed variance test

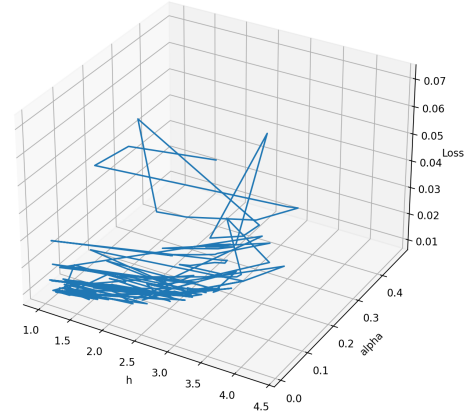


Figure 7: BBO of the GEM model with the human activity test

### 3.2 Principal Component Analysis

In the table 2 there can be seen the proposed hyperparameters and the resulting best achieved loss for the three test cases. The resulting values are promising since the behaviour is fairly similar to the GEM based results and it is a validation of the goodness. We can note how the objective function returns a higher loss in the second test case compared to GEM. The following graphs clearly represent the descent of the

optimizer as the result get closer to the local minimum and the behaviour is again analogous to the optimizing activity for GEM.

It is worth noting how the PCA algorithm is significantly faster then GEM but loses some accuracy with respect of the loss function we defined.

An important observation is that the loss function results in a plateau when either algorithm is tested against the mean shift dataset.

<b>Test Case</b>	<b><math>\alpha</math></b>	<b><math>h</math></b>	<b>Loss</b>	<b>Time</b>
$\mathcal{N} \sim (0.5, 0.5)$	0.024	6.934	0.6	221s
$\mathcal{N} \sim (1, 1)$	0.08	8.758	0.6	266s
$\mathcal{N} \sim (0, 2)$	0.611	1.004	0.297	265s
<i>Human Activity</i>	0.887	1.127	0.515	18259s

Table 2: Optimization results for PCA

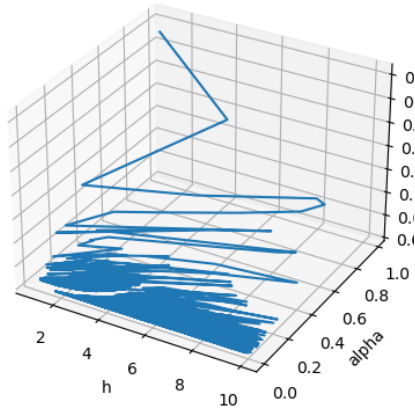


Figure 8: BBO of the PCA model with the subset test

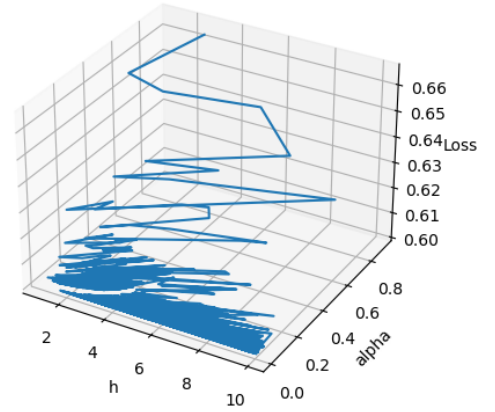


Figure 10: BBO of the PCA model with shifted mean test

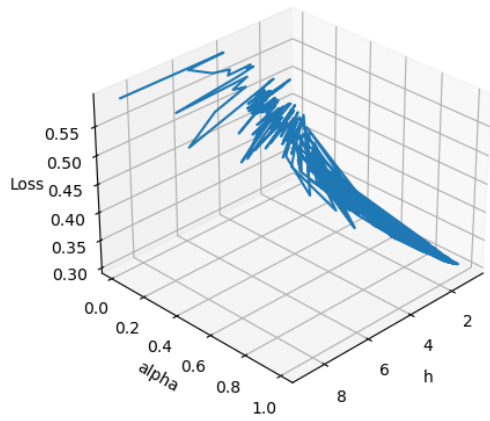


Figure 9: BBO of the PCA model with changed variance test

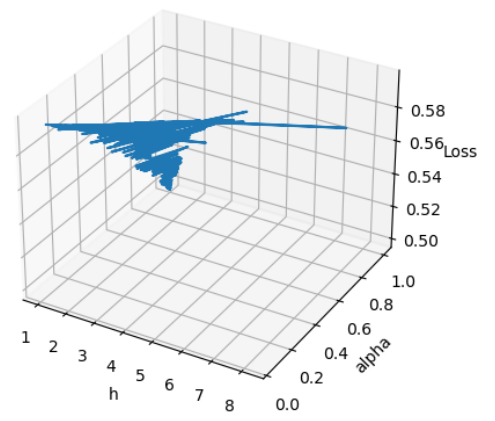


Figure 11: BBO of the PCA model with the human activity test

## 4 Testing

As an additional activity we performed some testing to compare the performance of the GEM based algorithm and the PCA based algorithm.

We gave both models the same hyper-parameters which are:

- $\alpha = 0.4$
- $h = 5.0$

And tested them on a number of different cases.

### 4.1 Increasing Dimension

In this test we plot both models' accuracy as the data dimension ( $p$ ) increases. We tested the models in two ways each:

**i.i.d.** the training set as a multivariate of independent identically distributed entries sampled from  $\mathcal{N} \sim (0, 1)$  and the test set as a multivariate sampled from  $\mathcal{N} \sim (0, 2)$

**Random covariance matrix** the training set was a multivariate with mean 0 and a covariance matrix resulted from randomly perturbing the identity matrix and the test set consisted of data points from  $\mathcal{N} \sim (0, 1)$

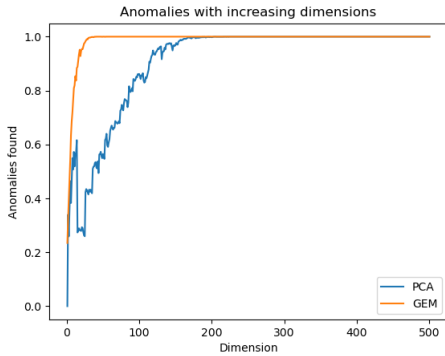


Figure 12: i.i.d. experiment

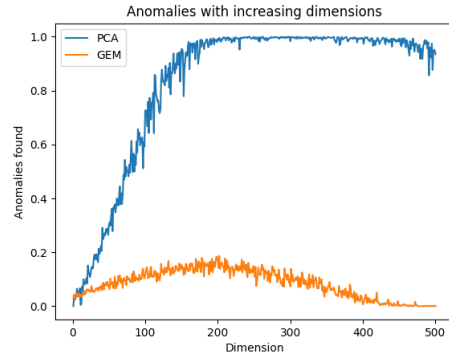


Figure 13: Random covariance experiment

As expected with the random covariance test PCA performed exponentially better than GEM because it manages to learn the relationships between the dimensions during the training phase while GEM only saw the dimensions as random and independent, so when tested against a random dataset of the same mean and variance GEM was not able to recognize the anomaly.

Using a multivariate where the dimensions were i.i.d. PCA was unable to extract meaningful correlations between them resulting in a performance similar yet slightly worse than GEM.

## 4.2 Increasing Offline Samples

In this test we plot the model’s accuracy as the training dataset grows against a fixed test dataset. We tested the models in two ways:

**i.i.d.** the training set as a multivariate of independent identically distributed entries sampled from  $\mathcal{N} \sim (0, 1)$  and the test set as a multivariate sampled from  $\mathcal{N} \sim (0, 2)$

**Random covariance matrix** the training set was a multivariate with mean 0 and a covariance matrix resulted from randomly perturbing the identity matrix and the test set consisted of data points from  $\mathcal{N} \sim (0, 1)$

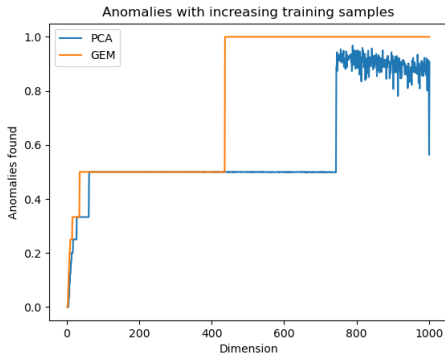


Figure 14: i.i.d. experiment

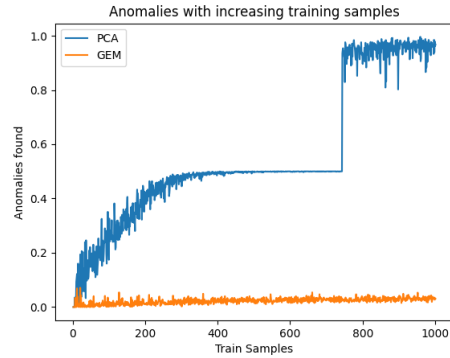


Figure 15: Random covariance experiment

Again, PCA performed better in the random covariance test and worse in the i.i.d. test. It is interesting to see how in the random covariance test the PCA-based algorithm reaches a plateau at 50% accuracy for and then after the training samples reached a certain number it skyrocketed to almost 100%. Also interesting to note how, in the i.i.d. test, the increase in training samples resulted in a series of constant values for the accuracy, with sudden steps.

## 4.3 Increasing Variance and Perturbation

In the first test we compare the accuracy of both algorithms with respect to a dataset in which the variance gradually increases. Let the training samples be from a normal multivariate random distribution with null mean and identity covariance matrix  $G_1 = I$ , we used as a testing dataset a normal multivariate distribution with  $\mu = 0$  and covariance matrix  $G_2 = kG_1$  with an increasing  $k = 1, \dots, 1000$ .

In the latter test we aim to show both models’ accuracy as the test dataset’s covariance matrix is gradually perturbed from the one they saw in the offline phase. We again used a multivariate normal random distribution with mean 0 and covariance matrix  $G_1 = I$  for our training while the testing used samples generated from a

multivariate normal random distribution with null mean but covariance matrix  $G_2 = I + r$  where  $r$  is noise.  $r$  was gradually increased by a factor of 0.1 for each iteration.

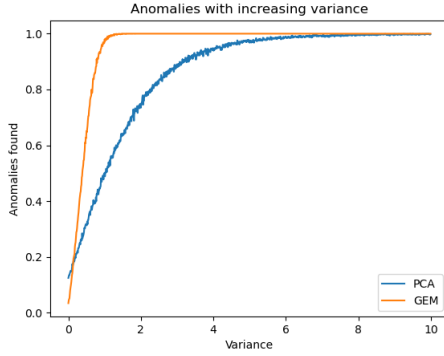


Figure 16: PCA and GEM performance with linear increase

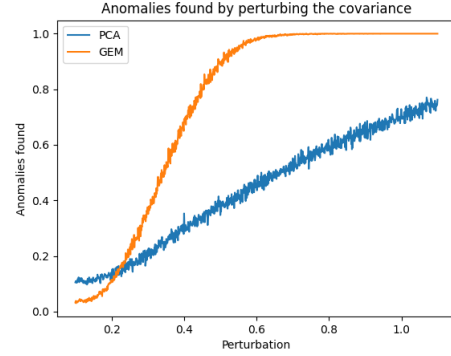


Figure 17: PCA and GEM performance with increasing noise in the covariance matrix

These test showed how a training set made of i.i.d. random variable always results in a worse performance for the PCA-based algorithm, as expected. The GEM-based algorithm, as showed in figure 16, is more sensible to changes in the variance of the samples.

In both cases the PCA approach is slower to converge compared to GEM. It is very visible in 17 in which the linear behaviour doesn't reach convergence by the end of the testing cycle.

## 5 Conclusions

We implemented an online anomaly detection system as proposed by [2] using C++ and Python and we

We subsequently optimized the models' two main parameters  $\alpha$  and  $h$  in order to infer the most suitable values on a number of different situations. Ulterior tests comparing the two models allowed us to observe the sensitivity and usefulness of the models in different test cases.

## 6 Future Developments

An important aspect to be considered in the algorithm is the unrealistic assumption that the data is comprised of independent and identically distributed data points (i.i.d.), to overcome this limitation it might be necessary to perform a detrending operation on the data in order to obtain a series of i.i.d. as assumed by the model.

Another improvement to be made is to fix the lack of reasoning on the history of the samples saw during training, as shown in ??.

There are also more hyper-parameters than  $\alpha$  and  $h$  to test and optimize, GEM also requires a  $k$  for the  $k$ NN phase and a fraction for the partition of the training set into the subsets  $S_1$  and  $S_2$ , finally PCA requires a  $\gamma$  for the desired accuracy of the algorithm.

## References

- [1] Srikanth Thudumu, Philip Branch, Jiong Jin, and Jugdutt Singh. A comprehensive survey of anomaly detection techniques for high dimensional big data. *Journal of Big Data*, 7, 07 2020.
- [2] Mehmet Necip Kurt, Yasin Ylmaz, and Xiaodong Wang. Real-time nonparametric anomaly detection in high-dimensional settings. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(7):2463–2479, July 2021.
- [3] J. Rapin and O. Teytaud. Nevergrad - A gradient-free optimization platform. <https://GitHub.com/FacebookResearch/Nevergrad>, 2018.
- [4] Jorge-L Reyes-Ortiz, Luca Oneto, Albert Samà Monsonís, Francesc Parra, and Davide Anguita. Transition-aware human activity recognition using smartphones. *Neurocomputing*, 171, 08 2015.