Image-based Process Monitoring via Adversarial Autoencoder with Applications to Rolling Defect Detection

Hao Yan, Huai-Ming Yeh, and Nurettin Sergin¹

Abstract-Image-based process monitoring has recently attracted increasing attention due to the advancement of the sensing technologies. However, existing process monitoring methods fail to fully utilize the spatial information of images due to their complex characteristics including the high-dimensionality and complex spatial structures. Recent advancements in unsupervised deep models such as generative adversarial networks (GAN) and adversarial autoencoders (AAE) has enabled to learn the complex spatial structures automatically. Inspired by this advancement, we propose an anomaly detection framework based on the AAE for unsupervised anomaly detection for images. AAE combines the power of GAN with the variational autoencoder, which serves as a nonlinear dimension reduction technique. Based on this, we propose a monitoring statistic efficiently capturing the change of the data. The performance of the proposed AAE-based anomaly detection algorithm is validated through a simulation study and real case study for rolling defect detection.

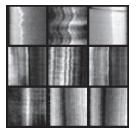
Index Terms—Statistical Process Control, Profile Monitoring, Deep Generative Models, Adversarial Autoencoder

I. INTRODUCTION AND LITERATURE REVIEW

Nowadays, image data are widely used in most manufacturing processes and service systems to evaluate the process performance and product quality due to the low implementation cost and the rich information it provides. Real-time image-based process monitoring and online product inspection are among the benefits that can be gained from this advancement.

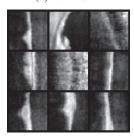
An example would be steel rolling process inspection where a high-speed video camera is set up to monitor the surface of the rolling bars. The camera captures images from the process as it unfolds. Then, certain areas of the images are focused on, which could potentially be coming from an out-of-control behaviour. Figure 1 shows examples these regions-of-interest, classified as normal or anomaly in the context of the manufacturing process. The main question we ask is: if we only have access to data that is produced by normal process, can we identify with a high accuracy, when the process is out-of-control.

The challenging part of this problem is that we're dealing with image data. First, image data are high-dimensional, which is susceptible to most problems related to dimensionality. Second, a corpus of images can have quite complex within sample and sample-to-sample variation that is often very hard to capture using simple models.



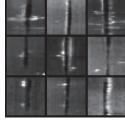


(a) Normal





(b) Anomaly 1



(c) Anomaly 2

(d) Anomaly 3

Fig. 1: Normal and abnormal rolling images

Most of the literature on process monitoring techniques for image data can be divided into the following two categories. Both categories attack the problems we mentioned above in different ways. Dimensionality reduction techniques find lower dimensional manifolds where it's easier to model the relationships within the darta. Previous attempts include classical linear dimensionality reduction techniques such as principal component analysis (PCA) [17] and independent component analysis [4]. Furthermore, some variants of the dimension reduction techniques are developed including functional PCA methods [9], [19], and tensor decomposition methods [14], [8]). The major drawback of these methods that they are constrained to learn linear manifolds which usually doesn't hold for image corpuses. A detailed analysis of this problem can be found in Shi et al. [13]. Consequently, the efforts move towards nonlinear dimensionality techniques such as kernel PCA [13] and maximum variance unfolding projections [12]. The referenced works focus exclusively on Phase-I monitoring therefore how to extend the methods to Phase-II is still an open problem. The other category of approaches are based on functional data analysis. These methods treat image as continuous functions and analyze the functional features by transforming the data into certain feature spaces for feature extraction, such as wavelet trans-

^{*} Hao Yan, Huai-Ming Yeh, and Nurettin Sergin are with School of Computing, Informatics, and Decision Systems Engineering, Arizona State University, Tempe, AZ 85281, USA.

This work is funded by NSF DMS-183036. Please address all correspondence to Professor Hao Yan, HaoYan@asu.edu

formation [7], [20], B-spline approximation [2], and functional decomposition-based techniques [15]. Finally, process monitoring techniques focusing on monitoring the feature coefficients or residuals [16]. These techniques typically have very strong assumptions on the image to be applied. For example, B-spline and kernel methods are designed for images with smooth background and wavelet is only designed for wave-form types of spatial structures. Both cases cannot be applied to more complex variation patterns.

Recently, unsupervised deep learning methods such as a generative adversarial network (GAN) [5] has been proposed and demonstrated that it can generate very realistic images. GAN-based approaches can capture the complex spatial correlation of the images data recently, therefore has been applied for anomaly detection in medical imaging [11], [18], [1]. However, these methods lack the ability to directly map the data into the feature space, which hinder its use for an efficient and real-time procedure. Recently, generative adversarial autoencoder (AAE) has been proposed [6], which combines the power of GAN and the variational autoencoder to ensure that the encoded features follow the normal distribution and the features can be directly computed by encoders. In this paper, we will propose how to use AAE methods for efficient and real-time process monitoring.

II. METHODOLOGY DEVELOPMENT

Before we dive into the proposed methodology for process monitoring, we would like to first briefly review the generative adversarial network (GAN), generative adversarial autoencoder (AAE), and how it can be used for efficient anomaly detection and process monitoring.

A. Review of Generative Adversarial Network

Generative adversarial networks (GAN) [5] has recently gained much attention due to its ability to learn the complex high-dimensional distribution by jointly train a generator G and discriminator D in a zero-sum game. The discriminator D(x) is normally a trained neural network that computes the probability that a point x is from the generator G(z) or the set of real samples. The generator uses a function G(z), which is typically another neural network, to generate the samples from the prior p(z) (i.e. usually a multivariate normal distribution with identity covariance matrix) to generate the data. The goal is to maximally confuse the discriminator D(x) so that it cannot tell the differences between whether the data is generated from the generator G(z) or from the real samples. In this case, the generator G(z) is identical to the real data distribution p(x).

The following minimax problem function is used to train D and G simultaneously [5].

$$\min_{G} \max_{D} E_{x \sim p_d(x)} [\log D(x)] + E_{z \sim p(z)} [\log(1 - D(G(z)))]$$
(1)

Typically, the training of GAN happens in two stages: (a) train the discriminator D to distinguish the true samples and the fake samples. (b) train the generator G to fool the discriminator.

However, the major limitation of GAN is that it lacks an efficient encoder to match the original data back to its feature. We will discuss how generative adversarial autoencoder (AAE) is able to solve this problem and how it can be used for anomaly detection in the next subsection.

B. Adversarial Autoencoder and its Application to Anomaly Detection

1) Adversarial Autoencoder: In the Phase I analysis, we can collect a large number of in-control (or normal) samples to learn the complex high-dimension distributions, where generative adversarial autoencoder can be used.

Here, we will still use x to represent the high-dimensional data to be monitored (e.g. signals, images, functional data from different sensors) and z to be the low-dimensional features. AAE is designed by introducing another autoencoder function q(z|x) by taking the input data x to compute the latent code z. The goal of AAE is to match the aggregated posterior distribution $q(z) = \int_x q(z|x)p_d(x)dx$ to a certain prior distribution p(z), where $p_d(x)$ is the data distribution. This can also act like a regularization so that the method would not over-fit. In order to do so, an adversarial network is trained to guide the q(z) to match the p(z) which is the same as how the adversarial training in GAN is used to match the generated data distribution to the sample data distribution.

In the original papers, the authors propose three different types of encoders for q(z|x). In this paper, we will just use the deterministic posterior where the q(z|x) is assumed to be a deterministic function of x since we have not found clear differences in terms of performance in other types of encoders.

Furthermore, in the AAE method, adversarial training is used to ensure the encoded features can match the Gaussian prior distribution by introducing the discriminator.

$$\min_{Q} \min_{G} l_{recon} + \min_{Q} \max_{D} l_{adversarial}$$
 (2)

Finally, the loss function consists of two components: i) The reconstruction error: The cost function of the reconstruction error is l_{recon} . Here, the reconstruction error l_{recon} is defined as

$$l_{recon} = E_{x \sim n_d(x)} l(x, G(Q(z|x))) \tag{3}$$

It's worth noting that the data likelihood is derived from the distribution of the data. ii) The discrimination error: The discriminative error is designed based on the adversarial training, which leads to the loss function $l_{adversarial}$ similar to (1) to match the p(z) and $q(z) = \int_x q(z|x)p_d(x)dx$. Here the adversarial loss $l_{adversarial}$ is defined as:

$$l_{adversarial} = E_{z \sim p(z)}[\log D(z)]$$

$$+ E_{x \sim p_d(x)} E_{z \sim Q(z|x)}[\log(1 - D(z))]$$
(4)

The adversarial loss $l_{adversarial}$ is to ensure that the latent code would match the aggregated posterior distribution and prior distribution p(z). After the generator G and the discriminator D are trained, we would like to discuss the way to monitor the data in the next subsection.

2) Phase I and Phase II Analysis: We propose to use the reconstruction error $l_{recon}(x) = l(x, G(Q(z|x)))$ as the monitoring statistics. Here, l(x, G(Q(z|x))) is the likelihood function, which depends on the probability distribution of the original data x. For example, if the data is Gaussian distributed, we can use the sum of squared error between the input image and the reconstruction image $\|x - G(Q(z|x))\|^2$. If the data is Bernoulli distributed, the cross entropy loss can be used, which is defined as $\|x \log(G(Q(z|x))) + (1 - x) \log(1 - G(Q(z|x)))\|_1$.

In the phase I analysis, we would like to sample a set of validation data set which belongs to the normal (in-control) samples but not in the training data set to estimate the distribution of the $l_{recon}(x)$ and the control limit c_0 . If the number of the normal samples is large, using the quantile is normally good enough of one validation data set is often good enough. However, if the number of the normal samples is limited, we can also use the cross-validation to compute the quantile of multiple validation data set as the control limit c_0 . Finally, c_0 is often selected such that the type I error is as a certain number. In this paper, we will set up the control limit c_0 based on the 5% false positive rate. Finally, in the Phase II analysis, we propose to use the $l_{recon}(x) > c_0$ to detect the anomalies.

III. SIMULATION STUDY

A. Simulation Setup

In this section, we generate 32×32 2D images containing a circle with different locations and shapes for normal samples from the following formula:

$$y = \sqrt{1 - ((x_1 - x_{10})/a)^2 - ((x_2 - x_{20})/b)^2},$$

where $x_{10}, x_{20} \sim N(0.5, 0.05^2)$, $a,b \sim N(0.2, 0.05^2)$ denote pixel locations on an image.

To validate the Phase II process monitoring performance, we generate four types of abnormal data for performance evaluation by the following formula:

· Mean Shift:

$$y = \delta + \sqrt{1 - ((x_1 - x_{10})/a)^2 - ((x_2 - x_{20})/b)^2},$$

• Magnitude Change:

$$y = (1+\delta)\sqrt{1 - ((x_1 - x_{10})/a)^2 - ((x_2 - x_{20})/b)^2}$$

• Width Change:

$$y = \sqrt{1 - ((x_1 - x_{10})/(a + \delta))^2 - ((x_2 - x_{20})/b)^2},$$

• Location Change:

$$y = \sqrt{1 - ((x_1 - x_{10})/a)^2 - ((x_2 - x_{20})/(b + \delta))^2}.$$

Here δ represent the change magnitude, x_{10}, x_{20} denote pixel locations on an image, and a, b denote the widths of the circle. Examples of these generated normal images and abnormal images ($\delta = 0.3$) are shown in Figure 2.

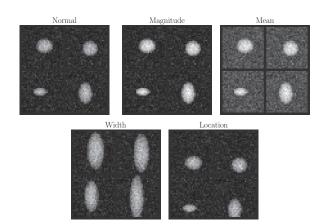


Fig. 2: Normal and abnormal images with $\delta=0.3$ (The change magnitude is enlarged to show a more clear change pattern)

B. Results

For the proposed AAE anomaly detection method, we proposed to use the latest DCGAN architectures [10] for the encoder and decoder, which combined a set of convolutional layers, batch normalization, and ReLU activation. We will compare the different latent dimensions such as 6 and 10 for the AAE method. Furthermore, in order to balance the generator and discriminator, the generator learning rate is set to 0.002 whereas the discriminator learning rate is 0.0002. Both use Adam optimizer for the best result. For the benchmark methods, we propose to compare with the widely used PCA-based methods with different principal components (PCs). For AAE, we will use the reconstruction error as the monitoring statistics. For PCA methods, we propose to combine the Q-chart and T2 chart for better performance [3]. Furthermore, to identify the anomaly, we would like to use the 5% false positive rate for all methods.

For the evaluation, we propose to use the detection power of Phase II analysis, defined as the percentage of samples that are correctly detected as anomalies by the algorithms over the number of anomalous samples. Finally, the detection power of AAE and PCA with different encoding dimensions and PC dimensions are shown in Table I, respectively. From Table I, we can conclude that AAE is more powerful than PCA in all cases that we tested. For examples, in the non-linear change such as width and location change, AAE has a very clear advantage over PCA since PCA is not able to represent these non-linear change patterns. In the linear change such as mean-shift, AAE still has some advantage over PCA. We can also conclude that simply adding the number of PCs only helps on the linear change patterns such as mean shift, but doesn't help too much on the complex nonlinear change patterns such as the width change and location change.

Furthermore, adding the number of PCs from 50 to 200 may also reduce the performance of PCA methods on the Width change (i.e. decrease from 18.78 to 18.35). However, AAE is much more robust to the selection of the latent dimensions. For example, the performance of AAE with 6

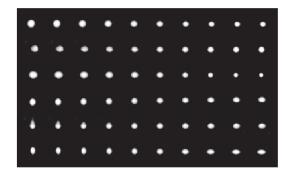


Fig. 3: A 2D Slice of Learned Manifold in the AAE model

or 10 latent dimensions are similar, and AAE with 6 latent dimensions works slightly better.

Furthermore, we also plot the detection power with different change magnitudes δ in Figure 4. From Figure 4, we can observe that PCA is sensitive to the number of PCs chosen. Magnitude change is challenging to both models since it doesn't change the shape of the original images. We found PCA can hardly detect the change even when $\delta=0.3$ whereas, AAE is able to reach about 100% accuracy. For the nonlinear changes such location change and width change, using AAE with 6 dimensions outperforms PCA with even up to 200 dimensions.

To understand how AAE detect the change, we would like to plot the real and reconstructed images in Figure 5. From this figure, we can conclude that AAE is able to generate accurate reconstruction images for the normal samples in Figure 5 (a). PCA with 10 PCs (i.e. PCA 10) cannot get the clear reconstruction of the original images. Furthermore, PCA with 200 PCs (i.e. PCA 200) can reconstruct the original images, but it generates a much noisy result. For the mean shift in Figure 5 (b), AAE is able to recover the darker circle and background. The reconstructed images of PCA with 10 and 200 PCs is largely affected by the background mean shift. For the magnitude change in Figure 5 (c), the AAE generates much darker reconstruction images than PCA, which shows that AAE actually learned the right magnitude for the normal samples, compared to the PCA methods. Finally, for width change and location change, the AAE reconstructed images are not affected by the out-ofcontrol samples, which make the reconstruction error a great monitoring statistics. However, PCA with 10 PCs generates very blurry images and PCA with 200 PCs tend to reconstruct the original images, which makes the Q-chart (i.e. control chard designed based on residual) not as effective.

Finally, to understand how AAE represent the complex nonlinear variational patterns, we also plot the 2D manifold reconstructed from the 2D lattice from a slice of 6-D latent space z in Figure 3. It is clear that the AAE actually learns a very smooth manifold and the upper-left corner has a much larger circle than the other corners. This demonstrates the smoothness of the generated samples from the latent feature

Methods		PC	AAE			
Dimension	6	10	50	200	6	10
Mean	18.30	22.54	96.90	97.66	100.0	100.0
Magnitude	37.07	34.85	32.77	32.31	97.97	97.47
Width	19.54	19.75	18.78	18.53	96.44	90.98
Location	19.72	19.95	19.51	19.77	96.58	91.31

TABLE I: Mean Shift Accuracy (%) Under $\delta = 0.06$

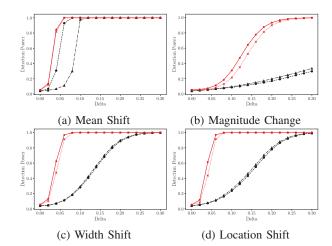


Fig. 4: Sensitivity Analysis of PCA and AAE

IV. CASE STUDY

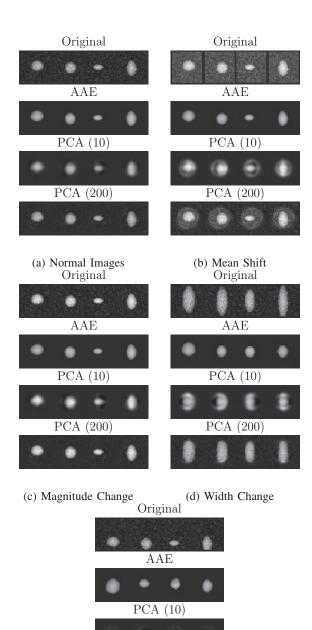
In this section, we will use real images from the quality inspection in the rolling manufacturing to illustrate the performance of the proposed AAE anomaly detection procedure. The dataset is made of metal rolling inspection images that are potentially defect. The domain engineers have labeled the images as normal or abnormal samples. Training data is made of 879 normal images and 294 abnormal images in 3 different types of abnormal conditions: Anomaly 1, Anomaly 2, and Anomaly 3. The regular images feature vertical texture or some minor overfills with black or white lines. The examples of normal images, Anomaly 1, Anomaly 2, and Anomaly 3 are shown in Figure 6 (a), (e), (i), (m), respectively.

We will compare AAE with PCA on the detection power (percentage of detected samples) with a fixed 5% false positive rate. We will use the same architecture and optimizer for the case study. We will use 6 or 10 latent space dimensions for the AAE method. For PCA, we will investigate to use 2,10,50,200 PCs for anomaly detection.

As shown in table II, PCA with 2 PCs works the best for *Anomaly 1*. The reason is that the change of *Anomaly 1* is quite obvious and the shape looks completely different than the normal vertical texture and can be fully detected via PCA with 6 PCs and AAE with 6 latent dimensions. However, for *Anomaly 2* and *Anomaly 3*, both are much harder to detect since the anomaly patterns also feature vertical textures. For AAE, in Anomaly 2 and 3, increasing the latent dimension to 10 will increase the detection accuracy to 100%. This is because AAE put the adversarial regularization so it can control the model complexity. However, for *Anomaly 2*, PCA with 50 PCs is optimal and achieves detection accuracy

TABLE II: Detection Power (%)

Methods		PC	CA	AAE			
Dim	2	6	50	200	2	6	10
Change 1	97.76	100.0	99.25	99.25	94.02	100	100
Change 2	34.72	43.05	59.72	47.22	54.16	86.11	100
Change 3	41.57	20.22	20.22	23.59	51.68	78.65	100



(e) Location Change

PCA (200)

Fig. 5: Original and reconstructed images for AAE and PCA with encoding dimensions 10 and 200

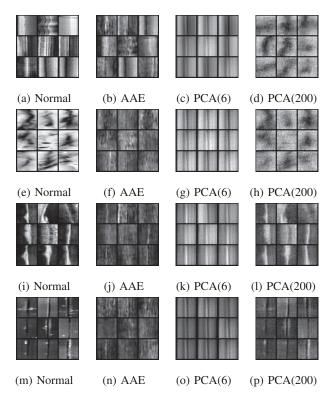


Fig. 6: Normal and Anomaly images, PCA reconstructed images (with 6 and 200 PCs), and AAE reconstructed images (with 10 latent dimension)

around 60%. For *Anomaly 3*, PCA with 2 PCs works the best. Further increasing the number of PCs decreases the detection power. This shows that simply increasing the number of PCs may lead to over-fitting problems.

Finally, we also plot the normal images, anomaly images, AAE-reconstructed images, and PCA reconstructed images with 6 PCs (i.e. PCA 6) and 200 PCs (i.e. PCA 200) in Figure 6. From Figure 6, we can conclude that AAE can reconstruct very realistic images with texture details. PCA with 6 PCs tends to create blurry images without enough details. On the other hand, PCA with 200 PCs will overfit and generate exactly the original images.

Finally, we also plot the AAE learned manifold in Figure 7. This figure is generated in a 2D lattice slice of a 10-D latent space z by the decoder q(z|x). This map shows that AAE actually learns a smooth manifold of all the rolling images.

V. CONCLUSION

Image-based process monitoring and anomaly detection often deal with complex nonlinear spatial correlation structures. In this paper, we propose a nonlinear anomaly detection algorithm based on generative adversarial autoencoder. The proposed method has shown a largely improved accuracy

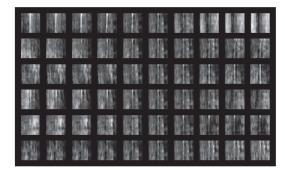


Fig. 7: Learned 2D Manifold

over the traditional linear process monitoring method such as PCA. A simulation study and a real case study from rolling manufacturing have been added to the original paper to demonstrate the advantage of the proposed AAE-based anomaly detection methods.

REFERENCES

- Samet Akcay, Amir Atapour-Abarghouei, and Toby P Breckon. Ganomaly: Semi-supervised anomaly detection via adversarial training. arXiv preprint arXiv:1805.06725, 2018.
- [2] Shing I Chang and Srikanth Yadama. Statistical process control for monitoring non-linear profiles using wavelet filtering and bspline approximation. *International Journal of Production Research*, 48(4):1049–1068, 2010.
- [3] Q Chen, U Kruger, M Meronk, and AYT Leung. Synthesis of t2 and q statistics for process monitoring. *Control Engineering Practice*, 12(6):745–755, 2004.
- [4] Yu Ding, Li Zeng, and Shiyu Zhou. Phase i analysis for monitoring nonlinear profiles in manufacturing processes. *Journal of Quality Technology*, 38(3):199, 2006.
- [5] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Advances in neural information processing systems, pages 2672–2680, 2014.
- [6] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. arXiv preprint arXiv:1511.05644, 2015.
- [7] Kamran Paynabar and Jionghua Jin. Characterization of non-linear profiles variations using mixed-effect models and wavelets. *IIE Transactions*, 43(4):275–290, 2011.
- [8] Kamran Paynabar, Jionghua Jin, and Massimo Pacella. Monitoring and diagnosis of multichannel nonlinear profile variations using uncorrelated multilinear principal component analysis. *IIE Transactions*, 45(11):1235–1247, 2013.
- [9] Kamran Paynabar, Changliang Zou, and Peihua Qiu. A change-point approach for phase-i analysis in multivariate profile monitoring and diagnosis. *Technometrics*, 58(2):191–204, 2016.
- [10] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. arXiv preprint arXiv:1511.06434, 2015.
- [11] Thomas Schlegl, Philipp Seeböck, Sebastian M Waldstein, Ursula Schmidt-Erfurth, and Georg Langs. Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In International Conference on Information Processing in Medical Imaging, pages 146–157. Springer, 2017.
- [12] Ji-Dong Shao and Gang Rong. Nonlinear process monitoring based on maximum variance unfolding projections. *Expert Systems with Applications*, 36(8):11332–11340, 2009.
- [13] Zhenyu Shi, Daniel W Apley, and George C Runger. Discovering the nature of variation in nonlinear profile data. *Technometrics*, 58(3):371– 382, 2016.
- [14] Hao Yan, Kamran Paynabar, and Jianjun Shi. Image-based process monitoring using low-rank tensor decomposition. *IEEE Transactions* on Automation Science and Engineering, 12(1):216–227, 2015.

- [15] Hao Yan, Kamran Paynabar, and Jianjun Shi. Anomaly detection in images with smooth background via smooth-sparse decomposition. *Technometrics*, 59(1):102–114, 2017.
- [16] Hao Yan, Kamran Paynabar, and Jianjun Shi. Real-time monitoring of high-dimensional functional data streams via spatio-temporal smooth sparse decomposition. *Technometrics*, (just-accepted), 2017.
- [17] Guan Yu, Changliang Zou, and Zhaojun Wang. Outlier detection in functional observations with applications to profile monitoring. *Technometrics*, 54(3):308–318, 2012.
- [18] Houssam Zenati, Chuan Sheng Foo, Bruno Lecouat, Gaurav Manek, and Vijay Ramaseshan Chandrasekhar. Efficient gan-based anomaly detection. arXiv preprint arXiv:1802.06222, 2018.
- [19] Chen Zhang, Hao Yan, Seungho Lee, and Jianjun Shi. Weakly correlated profile monitoring based on sparse multi-channel functional principal component analysis. *IISE Transactions*, 50(10):878–891, 2018
- [20] Shiyu Zhou, Baocheng Sun, and Jianjun Shi. An spc monitoring system for cycle-based waveform signals using haar transform. *IEEE Transactions on Automation Science and Engineering*, 3(1):60–72, 2006