# Autoencoder-based Representation Learning from Heterogeneous Multivariate Time Series Data of Mechatronic Systems

Karl-Philipp Kortmann, M. Sc.[†‡], Moritz Fehsenfeld, M. Sc.[†] and Dr.-Ing. Mark Wielitzka[†]

[†] Leibniz University Hannover, Institute of Mechatronic Systems, An der Universität 1, 30823 Garbsen, Germany
[‡] kortmann@imes.uni-hannover.de

## Abstract

Sensor and control data of modern mechatronic systems are often available as heterogeneous time series with different sampling rates and value ranges. Suitable classification and regression methods from the field of supervised machine learning already exist for predictive tasks, for example in the context of condition monitoring, but their performance scales strongly with the number of labeled training data. Their provision is often associated with high effort in the form of person-hours or additional sensors. In this paper, we present a method for unsupervised feature extraction using autoencoder networks that specifically addresses the heterogeneous nature of the database and reduces the amount of labeled training data required compared to existing methods. Three public datasets of mechatronic systems from different application domains are used to validate the results.

## 1   Introduction

Modern mechatronic systems contain a large number of internal sensor and control signals that can be accessed for example via fieldbus interfaces. In addition, these systems can be extended by external measuring systems, which in total leads to a heterogeneous set of multivariate time series signals (MTS), where heterogeneous here implies divergent sampling times, measuring resolutions or scale levels of the individual univariate signals.

In the course of progressive digitization, such system signals are increasingly used for classification tasks (such as *alarm* or *condition monitoring*) or for the regression of target variables that cannot be measured directly (e.g. process stability or quality). In the following, these are collectively referred to as *estimation tasks*.

The estimation methods primarily used for this purpose from the field of statistical or machine learning depend on extensive feature extraction [1], especially in the case of a high-dimensional, heterogeneous data basis. In the case of manual feature extraction, this requires a high level of technical expertise regarding the system at hand and is consequently accompanied by a high effort during implementation. Alternatively, supervised end-to-end estimation methods from the field of deep learning explicitly manage without previously extracted features, but require a large amount of already labeled data for training [2].

Methods of *representation learning* have recently emerged in the field of computer vision and speech recognition, enabling unsupervised feature extraction that outperforms the prediction performance of common manual and automated feature extraction methods in the case of only a small amount of existing labeled data [3]. The application to mechatronic systems has so far been limited to specialized single solutions [4], which cannot be directly transferred to any arbitrary mechanical or electrical system.

## 2   State of Research

### 2.1   Representation Learning

The term representation learning or feature learning covers methods that allow an automatic extraction of relevant features and thus make the step of manual feature engineering superfluous in principle. In the context of this work, the focus is on an unsupervised method that learns a compact representation of the MTS without knowledge of the target variable and thus only on the basis of the input time series. Originally derived from the field of computer vision, feature learning methods are increasingly adapted for estimation tasks based on time series data such as Chen et al. [4] using an autoencoder for feature extraction from the torque signals of a 6-axis industrial robot to predict collisions in the workspace. Li et al. [5] and Jiang et al. [6] both use a *Generative Adverssarial Network* (*GAN*) for feature extraction from industrial time series data, while Franceschi et al. [3] use a pure encoder-based network in combination with a so-called *triplet loss function* for classification on various reference time series. In this work, we focus on an autoencoder-based approach because, in contrast to GANs for instance, these generally provide a better representation of the population of training data, often at the cost of worse performance than purely generative models [7] (i.e., in generating realistic new time series), but this is not the focus of this paper.

**Autoencoder**   An autoencoder is an artificial neural network whose primary goal is to reconstruct an input signal $\mathbf{x}$ (see Figure 1). The dimension-reduced latent variable (also *bottelneck* or *latent space*)

$$\mathbf{z} = \phi\left(\mathbf{x}; \theta_{En}\right) \tag{1}$$

is the result of an encoder function with the parameters (weights) $\theta_{\text{En}}$ and is used as a feature vector in the representation learning in this work. The latent variable subsequently passes through a decoder

$$\tilde{\mathbf{x}} = \psi(\mathbf{z}; \theta_{De}) \qquad (2)$$

with parameters $\theta_{\text{De}}$, which generates a reconstruction $\tilde{\mathbf{x}}$ of the input signal. In the present case of mostly real-valued
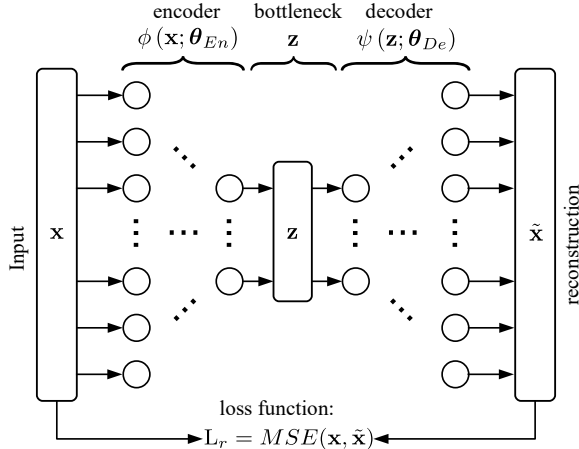


**Figure 1**  Schematic network diagram of a simple autoencoder with one-dimensional input and implied fully connected layers in encoder and decoder

input values, usually the mean squared error $L_{\text{r}} = \text{MSE}(\mathbf{x}, \tilde{\mathbf{x}})$ is used as reconstruction error and serves as a loss function during optimization.

Besides fully connected or dense layers with nonlinear activation functions, more complex neural layers also find use. For example, Bianchi et al. use recurrent (*RNN*) bidirectional layers in combination with an additional kernel loss function for feature extraction from MTS with missing values [8].

**Variational Autoencoder**  An extension of the regular autoencoder is the *Variational Autoencoder (VAE)*, which is mostly used as generative model in the field of computer vision. As an example of a *Variational Bayes* model, the VAE models the unknown distribution function of the input data $\mathbf{x} \sim p^*(\mathbf{x})$ using a model distribution $p_\theta(\mathbf{x}) \approx p^*(\mathbf{x})$ [9]. The stochastic decoder can therefore be understood as a conditional probability distribution $p_{\theta_{\text{De}}}(\mathbf{x}|\mathbf{z})$, which together with the prior distribution of the latent variable $p_\theta(\mathbf{z})$ forms a generative model by factorizing the multivariate distribution

$$p_\theta(\mathbf{x}, \mathbf{z}) = p_\theta(\mathbf{z}) p_{\theta_{\text{De}}}(\mathbf{x}|\mathbf{z}). \qquad (3)$$

Similarly, the encoder represents an inference model that can be conceived as a conditional probability distribution of the latent variables given input data $q_{\theta_{\text{En}}}(\mathbf{z}|\mathbf{x})$. This approach leads by the application of the *evidence lower bound, ELBO*

to the modified loss function [9]

$$\begin{aligned} L_{\theta_{\text{En}}, \theta_{\text{De}}}(\mathbf{x}) = {} & D_{\text{KL}}(q_{\theta_{\text{En}}}(\mathbf{z}|\mathbf{x}) \| p_\theta(\mathbf{z})) \\ & - \mathbb{E}_{q_{\theta_{\text{En}}}(\mathbf{z}|\mathbf{x})}\big(\log p_{\theta_{\text{De}}}(\mathbf{x}|\mathbf{z})\big), \quad (4) \end{aligned}$$

with the Kullback-Leibler divergence $D_{\text{KL}}$, which penalizes the deviation between a given prior distribution of the latent variable $\mathbf{z}$ and its actual (empirical) distribution given by the encoder. The second term represents the reconstruction error. The standard multivariate normal distribution $p_\theta(\mathbf{z}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is mostly used as prior for the latent variable. For the training using the standard stochastic gradient descent (*stochastic gradient descent, SGD*) method, the gradient of the loss function $\nabla L_{\theta_{\text{En}}}, \theta_{\text{De}}(\mathbf{x}_{\text{mb}})$ is calculated for each mini-batch of training data $\mathbf{x}_{\text{mb}}$ in order to perform minimization of the loss function as a function of the network parameters $\theta_{\text{En}}$ and $\theta_{\text{De}}$ (*backpropagation*). As can be seen in Figure 2a, this is not possible when directly drawing $\mathbf{z} \sim q_{\theta_{\text{En}}}(\mathbf{z}|\mathbf{x})$, since the backpropagation is interrupted by the random variable $\mathbf{z}$ [11]. Only with a mathematically equivalent reparametrization by swapping out the random variation into the random variable $\epsilon$ (which is typically modeled as normally distributed), a backpropagation of the error through the encoder is possible (so-called *reparameterization trick*, see Figure 2b).
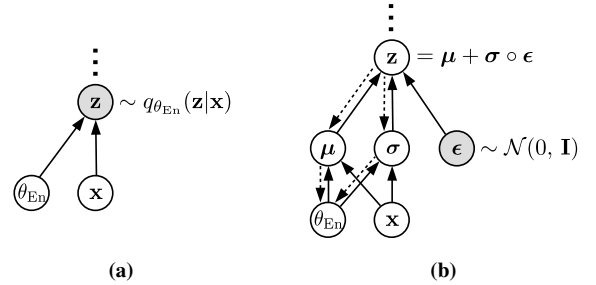


**Figure 2**  **(a)**: Direct sampling of $\mathbf{z}$ makes backpropagation ($--\rightarrow$) of the loss function impossible. **(b)**: Only with the reparametrization of $\mathbf{z}$ the optimization of the encoder parameters becomes possible. Independent random variables are shaded gray. Based on [10].

# 3  Methodology

This section presents the suggested autoencoder-based model, as well as the embedding pipeline. This is followed by details on the experimental design, in particular the chosen comparison methods and the selected publicly available datasets.

## 3.1  Feature Extraction

Unsupervised feature extraction from MTS with the amount of univariate signals $n_{\text{sig}}$ is performed using multiple VAEs trained separately for each univariate time series. This, in contrast to common two-dimensional convolutional network

architectures (*convolutional neural networks, CNN*), allows a parallelization of the training and an individual architecture of the networks, depending on the sampling rate of the heterogeneous signals[1].

As an aggregated feature, the 1D concatenation of the individual latent variables $\mathbf{x}_i$, $i \in \{1, ..., n_\text{sig}\}$ is used as input of the estimator. Table 1 lists the most important (hyper) parameters of the trained model. These were kept constant across all datasets in order to provide the best possible demonstration of generality.

Instead of specifying a constant dimension of the latent

**Table 1** Summary of the most important (hyper-)parameters of the univariate VAE model.

| (Hyper-)parameter | Value | Remark |
|---|---|---|
| Input dimension | $1 \times *$ | $*$: Maximum window length |
| Compression ratio $\kappa$ | 25 | Quotient of $*$ and $\dim(\mathbf{z})$ |
| Hidden layer | $[*/2, */2]$ | Two hidden layers |
| Activation radio | - | tanh resp. lin. for output layer |
| Regularization | - | Early stopping and $L_2$ norm |
| Optimizer | Adam | SGD optimizer [12] |
| Batch-size | $64 - 512$ | $\sim n_\text{train}$ |
| Normalization | - | Instance or layer norm. |
| Learning rate | $1 \times 10^{-4}$ | No learning rate scheduler |
| Max. epochs | $1 \times 10^3$ | Note: early stopping |

space, a constant compression rate $\kappa$ is chosen so that the number of extracted features scales proportional to the sampling rate and signal duration. Furthermore, during training and inference, *instance normalization* is performed for each windowed input time series to account for divergent ranges of signal values and stabilize convergence during the training.

## 3.2    Pipeline

The sequence of modeling/training and inference steps is typical for a so-called *semi-supervised* training procedure, consisting of unsupervised representation learning and supervised training of an estimator. (see Figure 3): After consecutive training of both models, they are applied unchanged during inference ( in this case on independent test datasets). The amount of labeled training data is varied during the testing procedure (logarithmic scaling of quantity).

The feature learning method (VAE) is provided with all available training data without labels beforehand in order to learn representations. This corresponds to the realistic use case of having a large amount of raw data, but only a certain fraction of it has been labeled. For each dataset, method, and quantity of labeled training data, $n_\text{repeat} = 10$ replicates are performed so that empirical mean standard deviation of the particular performance metric can be calculated. The train/test split of the datasets was provided by the authors of the same.

Python 3.8[2] on a computer with GPU support (CUDA 7.5) is used as experimental environment. The source code and
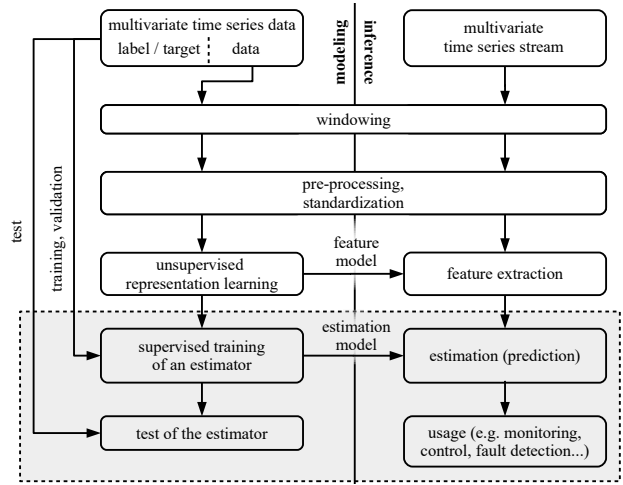
---



**Figure 3** Flowchart of the semi-supervised ML pipeline. The consecutive estimation task is highlighted in gray.

links to the used datasets are made publicly available[3].

## 3.3    Comparison methods

The presented representation learning model (VAE) as well as the semi-supervised pipeline are compared with a feature extraction method from the current state of the art (*Rocket*). Additionally, principal component analysis serves as a deterministic baseline comparison method. The Python library *tsfresh* allows automatic extraction and selection of statistical time and frequency domain features, thus providing a comparison method from the field of manual feature extraction.

For all comparison methods, a ridge regression-based estimator is used in that parameterization, as it is recommended as a favored predictor, especially by authors of several current feature extractors [13, 14]. In case of the VAE, a *support vector machine (SVM)* with Gaussian kernel is used to better account for the forced normal distribution character of the latent variable[4].

**PCA**    Using principal component analysis (PCA), time series data can be reduced in dimension by the help of singular value decomposition. PCA produces those orthogonal linear transformations of the input data which maximize the variance of the components of the target subspace in descending order. The obtained transformation can be used as a feature for classification or regression. In particular, PCA can be considered as a special case of a linear autoencoder without any hidden layers [9], which is why it will serve as a comparative baseline method here.

---

[1]On the other hand, existing cross-correlations between the individual univariate time series are no longer directly observable.

[2]In particular: torch 1.7.1, sktime 0.4.3 and sklearn 0.24.1.

[3]https://github.com/MrPr3ntice/vae_rep_learn_mts

[4]In particular, due to the symmetric kernel, closed intervals on one variable can be separated in case of classification.

**Table 2** Overview of selected data sets; $n_{\max}$ indicates the maximum number of data points in a time series sample (windowed).

| Dataset | Train samples | Test samples | Number of channels | Duration per sample ($n_{\mathbf{max}}$) | Target value (type) | Reference |
|---|---|---|---|---|---|---|
| Rolling bearing damages | 1440 | 800 | 7 | $0.2\,\mathrm{s}\,(800)$ | Rolling bearing condition (3 classes) | [18] |
| Stepper motors | 70152 | 23384 | 7 | $6\,\mathrm{ms}\,(60)$ | Operating condition (4 classes) | [19] |
| Hydraulic system | 1544 | 661 | 17 | $60\,\mathrm{s}\,(1200)$ | Cooler condition (3 classes), Hydraulic accumulator pressure (regr.) | [20] |

**Statistical Features**  The extraction of manually or automatically selected statistical features from time series for ML applications is widely used. The Python library *tsfresh* provides a collection of established extraction methods to automatically generate and select a variety of these features from time series data. The collection of features includes, for example, statistical ratios and correlations in both time and frequency domain. A complete overview of the extracted features can be taken from the libraries documentation [15]. In the present case, the default settings[5] is used.

**Rocket**  *Random convolutional kernel transform* is a state-of-the-art method that uses randomly sampled convolution kernels to extract features from time series. Subsequently, a Ridge regression[6] is trained with the thereby generated features and the known target values. Due to this straightforward setup, the computational cost of *Rocket* is lower than that of comparably performing methods [14]. To the state of this work, the method produces the best average classification results on the datasets of the UCR and UEA time series archive [17].

## 3.4    Data sets

The number of publicly available data sets is limited, especially in the area of mechanical and electronic systems. For validation of the proposed method, data sets covering a wide range of mechatronic applications are selected from those available. All data sets consist of real measurement data from several sensors, which differ e.g. in sampling rates. This results in three multivariate data sets from heterogeneous time series, from which three classification tasks and one regression task are derived. An overview of the selected data sets can be found in Table 2.

**Rolling bearing damages**  A widespread use case for machine learning applications is the detection of different rolling bearing damages. A comprehensive reference data set representing this use case has been published recently by Leissmeier et al. [18]. In addition to high-frequency sampled measurements of motor currents and housing vibration ($f_{\mathrm{s}} = 64\,\mathrm{kHz}$), additional, lower-frequency data such as radial force ($f_{\mathrm{s}} = 4\,\mathrm{kHz}$) and temperature ($f_{\mathrm{s}} = 1\,\mathrm{Hz}$) are available for

damage classification. The dataset includes different damage types of rolling bearings on the outer and inner ring as well as the data from undamaged bearings under different operating conditions. Further information on the data set can be found in [18].

**Stepper motors**  For stepper motor monitoring, there is a dataset published by Goubeaud et al. [19]. This includes measurements of current, voltage, and vibration (translational acceleration). The target variable is the operating mode of the stepper motor, which differentiates between clockwise and counterclockwise operation, as well as operation in the normal range and beyond the mechanical stop. A detailed description of the experimental setup and the data acquisition is given in the original publication [19].

**Hydraulic system**  The hydraulic system described by Helwig et al. [20] is equipped with a variety of different sensors. In addition to measurements of pressure and flow rates also temperature, current, and vibration are recorded.The sampling rates range from $f_{\mathrm{s}} = 1\,\mathrm{Hz}$ (temperature) to $f_{\mathrm{s}} = 20\,\mathrm{Hz}$[7] (pressure), resulting in a heterogeneous data set with a wide variety of sensor types, value range and sampling rate. In this work, the state of the cooler (fault classification) and the pressure in the hydraulic accumulator (regression) will be considered as target variables.

## 4    Results

The results are shown in Figure 4 (**a**)-(**d**). More detailed results can be retrieved from the Tables 3 and 4 in the appendix. For the first three cases, the *Rocket* comparison method achieves the highest test results with respect to the mean accuracy. Only in the case of the pressure estimation in the hydraulic accumulator, VAE shows almost consistently the lowest RMSE. In this regard, it should be noted that *Rocket* was originally developed for classification tasks and has been applied for this task in majority. Principal component analysis and the statistical features determined by *tsfresh* are not competitive for the prediction tasks on the first two data sets; for classification on the hydraulic data, all methods achieve similarly low error rates for a higher fraction of labeled training data.

A direct comparison with the partly available results of the

---

[5]extraction: `efficient`, selection: `extract_relevant_feat ures()`

[6]A special case of Tikhonov regularization for linear regression, sometimes also referred to as $L_2$-regularization [16]. It can be used in both classification and regression case.

[7]Compared to the original data set sampled at $100\,\mathrm{Hz}$, pressure has been sampled down to $20\,\mathrm{Hz}$ in favor of the computation time.
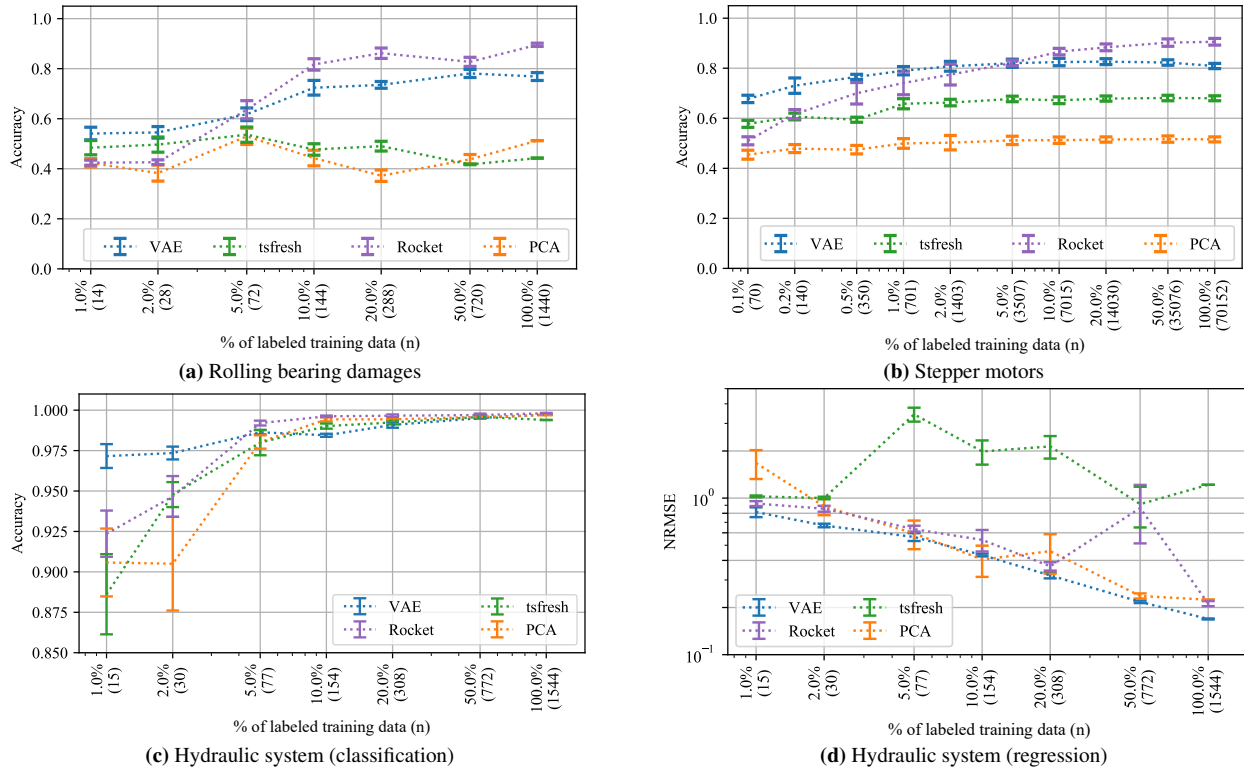
**Figure 4** Results of the four compared feature extractors on the three test datasets for varying proportions of labeled training data (log. scaling).The plots **(a) - (c)** show the accuracy for the respective classification tasks; **(d)** shows the normalized RMSE (log. scaling) for the regression task. For each proportion the empirical mean and standard deviation for $n_{\text{repeat}} = 10$ of the respective metric is shown.

original publications of the data sets is not possible at this point, since the multivariate case of a training data with a varying number of labeled data, which is treated here, was not considered in most studies. Additionally a neutral selection and parameterization of the compared estimators has been used in this work, in order to ensure the best possible comparability between the methods and not the highest possible performance of the same.

Looking at the results for a small number of labeled training data, it can be seen for all data sets that the presented implementation of VAE forms a better predictive measure based on the extracted features in the range up to at least 2 %. This supports the research hypothesis that autoencoder-based representation learning methods for extracting features from heterogeneous multivariate time series are particularly suitable for very small amounts of existing labeled training data in the presence of a sufficiently large amount of unlabeled training data at the same time.

## 5 Conclusion

In this work, an unsupervised autoencoder-based feature extractor for estimation tasks on heterogeneous, multivariate time series of mechatronic data was presented and validated on three data sets from the scientific community. Even though the prediction quality based on the entire training data

does not approach current state-of-the-art feature extractors like *Rocket*, an increased quality for the case of a low amount of labeled training data with a high availability of unlabeled data could be shown.

The presented results were obtained using only processed (reference) data sets. Thus, a consequential step towards automated estimation methods for the condition monitoring in comparable applications is the adaptation of the presented method for non-preprocessed raw data. Here, for example, the method of the *denoising autoencoder* may be considered as a possible extension for noisy or even corrupted input data.

**Publication notice** This is a pre-print version of the paper in German language submitted to *VDI Mechatronic Tagung 2021*, which has been published in the conference proceedings.

## References

[1] Fawaz, H. I.; Forestier, G.; Weber, J.; Idoumghar, L.; Muller, P. A.: *Deep learning for time series classification: a review*. Data Mining and Knowledge Discovery. (2019), 33(4), pp. 917–963.

[2] Lei, Y.; Jia, F.; Lin, J.; Xing, S.; Ding, S. X.: *An intelligent fault diagnosis method using unsupervised feature learning towards mechanical big data*. IEEE Transactions on Industrial

Electronics. (2016), 63(5), pp. 3137–3147.

[3] Franceschi, J.-Y.; Dieuleveut, A.; Jaggi, M.: *Unsupervised Scalable Representation Learning for Multivariate Time Series*. Advances in Neural Information Processing Systems. (2019) 32, ISSN 1049-5258.

[4] Chen, T.; Liu, X.; Xia, B.; Wang, W.; Lai, Y.: *Unsupervised Anomaly Detection of Industrial Robots Using Sliding-Window Convolutional Variational Autoencoder*. IEEE Access. (2020) 8, ISSN 2169-3536, pp. 47072–47081.

[5] Li, D.; Chen, D.; Jin, B.; Shi, L.; Goh, J.; Ng, S. K.: *MAD-GAN: Multivariate anomaly detection for time series data with generative adversarial networks*. In: International Conference on Artificial Neural Networks. (2019), Springer, Cham, pp. 703–716.

[6] Jiang, W.; Cheng, C.; Zhou, B.; Ma, G.; Yuan, Y.: *A novel GAN-based fault diagnosis approach for imbalanced industrial time series*. (2019), arXiv preprint arXiv:1904.00575.

[7] Grover, A.; Dhar, M.; Ermon, S.: *Flow-GAN: Combining maximum likelihood and adversarial learning in generative models*. In: Proceedings of the AAAI Conference on Artificial Intelligence. (2018) 32, 1.

[8] Bianchi, F. M.; Livi L.; Mikalsen, K. Ø.; Kampffmeyer M.; Jenssen, R.: *Learning representations of multivariate time series with missing data*. Pattern Recognition. (2019) 96, ISSN 0031-3203.

[9] Kingma, D. P.; Welling, M.: *An Introduction to Variational Autoencoders*. Foundations and Trends in Machine Learning. (2019) 12, No. 4, pp. 307–392.

[10] Kingma, D. P.; Welling, M.: *Auto-encoding Variational Bayes*. (2013), arXiv preprint arXiv:1312.6114.

[11] Rezende, D. J.; Mohamed, S., Wierstra, D.: *Stochastic back-propagation and approximate inference in deep generative models*. In: International Conference on Machine Learning. (2014), pp. 1278—1286.

[12] Kingma, D. P.; Ba, J.: *Adam: A method for stochastic optimization*. (2014), arXiv preprint arXiv:1412.6980.

[13] Le Nguyen, T.; Gsponer, S.; Ilie, I.; O'Reilly, M.; Ifrim, G.: *Interpretable time series classification using linear models and multi-resolution multi-domain symbolic representations*. Data Mining and Knowledge Discovery. (2019) 33, ISSN 1573-756X, pp. 1183–1222.

[14] Dempster, A.; Petitjean, F.; Webb, G. I.: *ROCKET: exceptionally fast and accurate time series classification using random convolutional kernels*. Data Mining and Knowledge Discovery. (2020) 34, ISSN 1573-756X, pp. 1454–1495.

[15] Christ, M.; Braun, N.; Neuffer, J.; Kempa-Liehr A.W.: *Time Series FeatuRe Extraction on basis of Scalable Hypothesis tests (tsfresh – A Python package)*. Neurocomputing. (2018) 307, ISSN 1573-756X, S. 72–77.

[16] Kennedy, P.: *A Guide to Econometrics*. 5th edition. Cambridge: The MIT Press. (2003). ISBN: 0-262-61183-X, pp. 205–206.

[17] Ruiz, A. P.; Flynn, M.; Large, J.; Middlehurst, M.; Bagnall, A.: *The great multivariate time series classification bake off: a review and experimental evaluation of recent algorithmic advances*. Data Mining and Knowledge Discovery. (2020), Springer OA, ISSN: 1573-756X, pp. 1–49.

[18] Lessmeier, C.; Kimotho, J.; Zimmer,D.; Sextro, W.: *Condition Monitoring of Bearing Damage in Electromechanical Drive Systems by Using Motor Current Signals of Electric Motors*. European Conference of the Prognostics and Health Management Society. (2016).

[19] Goubeaud, M.; Grunert, T.; Lützenkirchen, J.; Joußen, P.; Ghorban, F.; Kummert, A.: *Introducing a New Benchmarked Dataset for Mechanical Stop Detection of Stepper Motors*. 2020 27th IEEE International Conference on Electronics, Circuits and Systems (ICECS). (2020), pp. 1–4.

[20] Helwig, N.; Pignanelli, E.; Schütze, A.: *Condition monitoring of a complex hydraulic system using multivariate statistics*. 2015 IEEE International Instrumentation and Measurement Technology Conference (I2MTC) Proceedings. (2015) 24, ISSN, pp. 210–215.

# Appendix

**Table 3**  Test results of the classification tasks with all three data sets. The mean and standard deviation of the accuracy for different proportions of labeled training data are given in each case.

| % of labeled training samples ($n_{\text{train}}$) | VAE + SVM | PCA + Ridge | tsfresh + Ridge | Rocket + Ridge |
|---|---|---|---|---|
| Rolling bearing damages data set | | | | |
| 1.0% (14) | **.541 ± .023** | .424 ± .015 | .484 ± .026 | .424 ± .009 |
| 2.0% (28) | **.545 ± .021** | .383 ± .029 | .496 ± .012 | .426 ± .008 |
| 5.0% (72) | .618 ± .023 | .529 ± .029 | .536 ± .027 | **.636 ± .032** |
| 10.0% (144) | .724 ± .026 | .443 ± .028 | .477 ± .020 | **.817 ± .020** |
| 20.0% (288) | .735 ± .012 | .372 ± .020 | .490 ± .017 | **.862 ± .019** |
| 50.0% (720) | .781 ± .015 | .438 ± .017 | .417 ± .001 | **.827 ± .016** |
| 100.0% (1440) | .768 ± .014 | .512 ± .000 | .443 ± .000 | **.895 ± .006** |
| Stepper motors data set | | | | |
| 0.1% (70) | **.677 ± .013** | .454 ± .016 | .578 ± .013 | .510 ± .014 |
| 0.2% (140) | **.730 ± .025** | .479 ± .014 | .607 ± .012 | .616 ± .016 |
| 0.5% (350) | **.765 ± .010** | .474 ± .015 | .594 ± .010 | .700 ± .033 |
| 1.0% (701) | **.790 ± .015** | .499 ± .017 | .658 ± .017 | .741 ± .036 |
| 2.0% (1403) | **.808 ± .017** | .502 ± .023 | .663 ± .012 | .774 ± .032 |
| 5.0% (3507) | .820 ± .014 | .512 ± .015 | .677 ± .011 | **.821 ± .011** |
| 10.0% (7015) | .825 ± .013 | .512 ± .012 | .672 ± .011 | **.866 ± .012** |
| 20.0% (14030) | .826 ± .011 | .515 ± .010 | .679 ± .010 | **.883 ± .013** |
| 50.0% (35076) | .823 ± .011 | .517 ± .012 | .681 ± .011 | **.902 ± .013** |
| 100.0% (70152) | .809 ± .010 | .516 ± .010 | .680 ± .010 | **.906 ± .012** |
| Hydraulic system dataset | | | | |
| 1.0% (15) | **.972 ± .006** | .906 ± .018 | .886 ± .021 | .924 ± .012 |
| 2.0% (30) | **.974 ± .003** | .905 ± .025 | .948 ± .007 | .947 ± .011 |
| 5.0% (77) | .986 ± .001 | .980 ± .004 | .980 ± .007 | **.992 ± .001** |
| 10.0% (154) | .984 ± .001 | .994 ± .001 | .990 ± .001 | **.996 ± .000** |
| 20.0% (308) | .991 ± .002 | .994 ± .001 | .992 ± .001 | **.997 ± .001** |
| 50.0% (772) | .995 ± .000 | .996 ± .001 | .996 ± .001 | **.997 ± .001** |
| 100.0% (1544) | .997 ± .000 | .997 ± .000 | .994 ± .000 | **.998 ± .000** |

**Table 4**  Test results for the regression task with the hydraulic system data set. The mean value and standard deviation of the normalized RMSE for the prediction of the pressure in the hydraulic accumulator for different proportions of labeled training data are given in each case.

| % of labeled training samples ($n_{\text{train}}$) | VAE + SVR | PCA + Ridge | tsfresh + Ridge | Rocket + Ridge |
|---|---|---|---|---|
| 1.0% (15) | **.814 ± .058** | 1.674 ± .350 | 1.026 ± .013 | .923 ± .033 |
| 2.0% (30) | **.669 ± .018** | .885 ± .107 | 1.002 ± .019 | .856 ± .037 |
| 5.0% (77) | **.565 ± .033** | .595 ± .123 | 3.429 ± .350 | .633 ± .022 |
| 10.0% (154) | .433 ± .009 | **.405 ± .091** | 1.986 ± .350 | .541 ± .085 |
| 20.0% (308) | **.321 ± .014** | .458 ± .129 | 2.138 ± .350 | .367 ± .023 |
| 50.0% (772) | **.217 ± .004** | .237 ± .009 | .915 ± .266 | .864 ± .350 |
| 100.0% (1544) | **.169 ± .001** | .224 ± .000 | 1.217 ± .000 | .212 ± .008 |