



Data-driven prognostics with low-fidelity physical information for digital twin: physics-informed neural network

Seokgoo Kim^{1,2} · Joo-Ho Choi² · Nam Ho Kim¹

Received: 4 April 2022 / Revised: 26 July 2022 / Accepted: 27 July 2022 / Published online: 2 September 2022
© The Author(s), under exclusive licence to Springer-Verlag GmbH Germany, part of Springer Nature 2022

Abstract

In the absence of a high-fidelity physics-based prognostics model, data-driven prognostics methods are widely adopted. In practice, however, data-driven approaches often suffer from insufficient training data, which causes large training uncertainty that hinders the Digital twin (DT)-based decision-making. In such a case, the integration of low-fidelity physics with a data-driven method is highly demanded. This paper introduces physics-informed neural network (PINN)-based prognostics that can utilize low-fidelity physics information, such as monotonicity or the sign of curvature. Low-fidelity physics information is included as a constraint during the optimization process to reduce the training uncertainty in the neural network model by preventing unrealistic predictions. The proposed method is applied to two case studies to demonstrate the effect of reducing the prediction uncertainty and the robustness to the variability in test data. The two case studies show that PINN-based prognostics can successfully reduce the prediction uncertainty and yield more robust prognostics performance than the ordinary neural network.

Keywords Physics-informed neural network · Prognostics · Uncertainty quantification · Remaining useful life

1 Introduction

One of the promising technologies of industry 4.0 is a smart factory that aims to improve productivity in the manufacturing process. For the successful application of smart factories to the real industry, a digital twin (DT) which represents the virtual counterpart of a factory, manufacturing, and product is considered an essential ingredient since it supports decision-making throughout the lifecycle of the system. It is well known that the DT can be combined with prognostics and health management (PHM) to minimize downtime and

maintain reliable system operation. PHM is a key function for predictive maintenance to estimate the current health condition of the system and predict the remaining useful (Negri et al. 2021). PHM consists of two approaches, namely, physics-based and data-driven approaches (Lei et al. 2018; Kim et al. 2021; Lee et al. 2014). The former utilizes physics models describing the behavior or degradation of target systems, whereas the latter employs artificial intelligence (AI) or machine learning (ML) to build a surrogate model to replace the physical one. Two methods have their own pros and cons. The physics-based method is considered the most accurate approach. However, it is challenging to obtain or establish a high-fidelity physics model for a complex system. On the other hand, the data-driven approach does not require an in-depth understanding of physics. Based on the existing training dataset, the data-driven approach aims to find the hidden pattern behind the data. In the absence of physical knowledge, the data-driven approach is considered an alternative method to the physics-based approach. This is why most recent researches focus on data-driven methods for both fault diagnosis (Kim and Choi 2019; Kim et al. 2020a; Ham et al. 2019) and prognosis (Wang 2010; Heimes 2008). Moreover, the rise of deep-learning algorithms and the improvement in computational

Responsible Editor: Chao Hu

Topical Collection: Advanced Optimization Enabling Digital Twin Technology.

Guest Editors: C Hu, VA González, Z Hu, T Kim, O San, P Zheng.

✉ Nam Ho Kim
nkim@ufl.edu

¹ Department of Mechanical and Aerospace Engineering, University of Florida, Gainesville, FL, USA

² Department of Aerospace and Mechanical Engineering, Korea Aerospace University, Goyang-si, South Korea

power accelerate the influx of various machine learning algorithms into the PHM research field. There is no doubt that it is the renaissance of data-driven methods. Data-driven approaches, however, cannot be free from their inherent limitations. The main bottleneck of data-driven approaches is the requirement of a huge amount of training data. Insufficient training data can cause poor accuracy and large training uncertainty that hinders appropriate decision-making. In practice, it is almost impossible to obtain a large amount of fault data of a real mechanical or electrical system in the engineering field. This problem becomes worse for fault prognostics that aims to predict the remaining useful life of the engineering system. Different from fault diagnostics, prognostics focuses on the extrapolation region that cannot be covered with existing training data. Moreover, there are no decision-makers or maintenance engineers who operate their engineering assets until failure. This is considered the most significant challenge when adopting AI- or ML-based prognostics. Considering that the DT aims to support the decision-makers with a number of functions such as parameter optimization, monitoring, and behavior prediction based on the PHM (Lim et al. 2020), this challenge is directly pertinent to the performance of the DT.

To alleviate the drawback of data-driven approaches, integrating physical knowledge with the neural network (NN) is demanding. A physics-informed neural network (PINN) is an emerging approach that encodes physical laws or physical knowledge into the NN to benefit from physics-based and data-driven methods simultaneously. Most of their applications focus on solving governing partial differential equations by utilizing available data and boundary conditions (Raissi et al. 2019; Yang et al. 2020; Meng et al. 2020; Mao et al. 2020; Fang and Zhan 2019). Representatively, Raissi et al. (2019) presented the basic concept of PINN that integrates a governing equation, which is crucial mathematical information to describe the physical behavior of the system, into a NN framework. They showed a systematic methodology to solve the forward and inverse problems of a non-linear partial differential equation (PDE). Yang et al. (2020) demonstrated a new class of physics-informed generative adversarial networks to solve the forward, inverse, and mixed problems of PDE based on the limited number of measurements. Meng et al. (2020) pointed out that original PINN becomes computationally intractable for long-time integration of time-dependent PDE. To address this problem, they proposed parallel-in-time PINN, which solved two different PDEs using a coarse-grained solver.

The application of PINN is not limited to solving PDE but also applied to fault prognostics. Nascimento and Viana introduced a novel approach that encodes the physics model (i.e., Paris law) into a recurrent neural network (Nascimento and Viana 2019). They demonstrated that the physics-informed layer can be used to predict the next crack size.

Yucesan and Viana (2021, 2020, 2022) proposed a hybrid model that is designed to merge physics-informed and data-driven layers within the deep neural network to predict bearing fatigue. Nascimento et al. (2021) developed a hybrid modeling approach that incorporates reduced-order models and NN to build a battery prognostics model. Dourado and Viana (2020) presented PINN to compensate for the missing physics of corrosion in the fatigue model and proved that the prediction result from PINN was good. Goswami et al. (2020) proposed a PINN integrating phase-field modeling into a NN.

The state-of-the-art fault prognostics using PINN shows that researchers successfully put the first step of embedding physical models into the NN framework. Prior research on PINN-based prognostics, however, still has the following critical limitations that should be explored further:

- A high-fidelity physical model is still required to be embedded into the NN framework. As mentioned above, it is not trivial to obtain a high-fidelity physical model in the real industrial field.
- In reality, crude/abstract physical information is often available. In some cases, physical information is in the form of expert opinions and does not have a mathematical formula. Nevertheless, there is no work explicitly encoding the low-fidelity physical information into the NN framework.
- The existing method of modifying the loss function did not focus on the extrapolation region, which means that their modified loss function is only trained within the interpolation region.
- Existing research using the PINN requires the run-to-fail data to train their prognostics model. However, it is unrealistic to assume that there is a large number of training data in the real field.
- The performance of data-driven prognostics is evaluated for the fixed testing datasets. The training process of NN divides the existing training data into training, validation, and testing data. The uncertainty in the existing dataset and the arbitrary separation of training, validation, and testing data contribute to uncertainty in the NN model. The contribution of data uncertainty should be quantified appropriately.

Motivated by these limitations, we aim to propose a new framework of PINN with low-fidelity physical information for the DT. The performance of the proposed method is evaluated on a synthetic dataset generated from Paris' law that describes the crack growth. To quantify the uncertainty in testing data, multiple noise perturbations from the same distribution are added to the true crack growth behavior. The major contributions of this paper can be summarized as follows:

- A framework of physics-informed neural network-based prognostics is proposed with low-fidelity physical information
- Physical information in the extrapolation region is included in the form of a constraint during the optimization of neural work
- To address the problem of the lack of run-to-fail training data, the proposed method is developed based on the assumption that there is no run-to-fail data to train the neural network.
- Training uncertainty in NN is significantly reduced by integrating physical information
- Data uncertainty is quantified by using synthetic measurement data

The remainder of the paper is organized as follows. Section 2 describes the background of the physical information for prognostics and types of uncertainty in NN. Section 3 details the proposed method. In Sect. 4, the proposed PINN-based prognostics is applied to case studies for cooling fan bearing prognostics and a crack growth problem. For the cooling fan bearing prognostics, the effectiveness of the proposed method is demonstrated. Then, the robustness of the algorithm is investigated using the crack growth data generated from Paris' law. Finally, the paper is concluded with discussions and future works in Sect. 5.

2 Background

2.1 Physical information

There have been several trials to integrate physics-based and data-driven prognostics (Yucesan and Viana 2021, 2020, 2022; R. Giorgiani do Nascimento, F. Viana, M. Corbetta, and C. S. Kulkarni 2021; Dourado and Viana 2020; Goswami et al. 2020). In their approaches, they assumed that an accurate physical degradation or dynamic model exists. However, such a hybrid approach may not be effective if a high-fidelity physical model exists. Moreover, it is hard to establish a high-fidelity physical model that describes the degradation trends or dynamics of engineering systems. This problem becomes worse for system-level or multiple

component degradation cases. Therefore, it is necessary to define the fidelity of physical knowledge and develop an appropriate hybrid approach to maximize the utilization of physical information.

In this paper, physical knowledge is categorized into four levels as shown in Table 1, and a hybrid prognostics method is developed with low-fidelity physical information. For the purpose of explanation, each level of physical information is explained with a crack growth problem, where a_t , z_t , N , u , and θ represent the real crack size, measured crack size, cycle, other parameters related to operating condition, and model parameters, respectively. Level 0 refers to the situation when there is no available physical information about the target system. Then, the only possible choice is to build a data-driven model (i.e., a black-box model). Level 1 is the lowest level of available physical information that represents the crude behavior of parameters that is 'not accurate but informative.' For example, it is widely known that pressure decreases when altitude increases. In the case of crack propagation, the crack size can only increase over time; i.e., $da/dN \geq 0$. This kind of information sometimes depends on domain knowledge or human expertise. The second level is a low-fidelity model (i.e., empirical model) that can be easily found in the battery prognostics. Many engineering problems rely on the empirical model since it is difficult to develop a high-fidelity model that is based on complete physical reasoning. Finally, level 3 is considered the most informative physical information. At this level, a reliable physical model is available, which can describe the behavior of the system operation or degradation, such as Paris' law (Paris and Erdogan 1960) or Huang's model (Huang et al. 2008) for crack propagation. It is hard to justify utilizing a data-driven or hybrid method because the high-fidelity physical model performs always better. According to the current categorization, levels 0 and 3 are out of the discussion because it is obvious to use a data-driven method for level 0 and a physics-based method for level 3. The main concern of this paper is how to utilize the low-fidelity physical information (i.e., levels 1 and 2) in the context of data-driven prognostics.

Even if the main purpose of many hybrid methods is to improve the accuracy of prognostics, the focus of this paper is on uncertainty reduction. Figure 1 illustrates how

Table 1 Description of the levels of available physical information

Level	Description	Approach	Expression
0	No physical meaning	Data-driven	$z_t = f(z_{t-1}, u)$
1	The behavior of physical parameters Low-fidelity physical knowledge	-	$\frac{da}{dN} \geq 0$
2	Low-fidelity physical model Incomplete model	Empirical model	$z = g(N, \theta)$
3	High-fidelity model Physical model	Physics-based	$\frac{da}{dN} = C(\Delta K)^m$

the physical information contributes to improving the performance of the data-driven method. Under the limited number of training data, data-driven methods contain large uncertainty plotted with the blue line. Low-fidelity physical information about the range of output can define the feasible region of the output and shape the uncertainty in the prediction from the data-driven method. The red dashed line shows the reduced uncertainty by the physics-informed neural network.

2.2 Uncertainty in neural network

For a safe decision-making process, it is crucial to identify the sources of uncertainty and quantify them. Among many sources of uncertainty, this paper focuses on two uncertainties: training uncertainty and data uncertainty. In the following subsections, each type of uncertainty is introduced along with the uncertainty quantification method.

2.2.1 Training uncertainty

In the training process of NN, three factors contribute to uncertainty in NN as shown in Fig. 2. The first factor is related to the early stopping criterion, which is widely used to prevent overfitting. The early stopping criterion divides training data into training and validation sets and terminates the training process when the error in the validation set begins to increase. In this process, the model performance can be different depending on how the training and validation sets are selected. The second factor is associated with a network structure, including the number of nodes and layers and types of activation functions. This corresponds to a model-form uncertainty of NN. The last factor is due to optimization algorithms. Since optimization algorithms can only find local optima based on initial values of parameters, the trained parameters may be different with different starting points. The uncertainty caused by these three factors is referred to as training uncertainty in this paper.

To remove the uncertainty of selecting the validation set, this paper used a physics-based constraint term to prevent the model from overfitting. This approach makes it possible

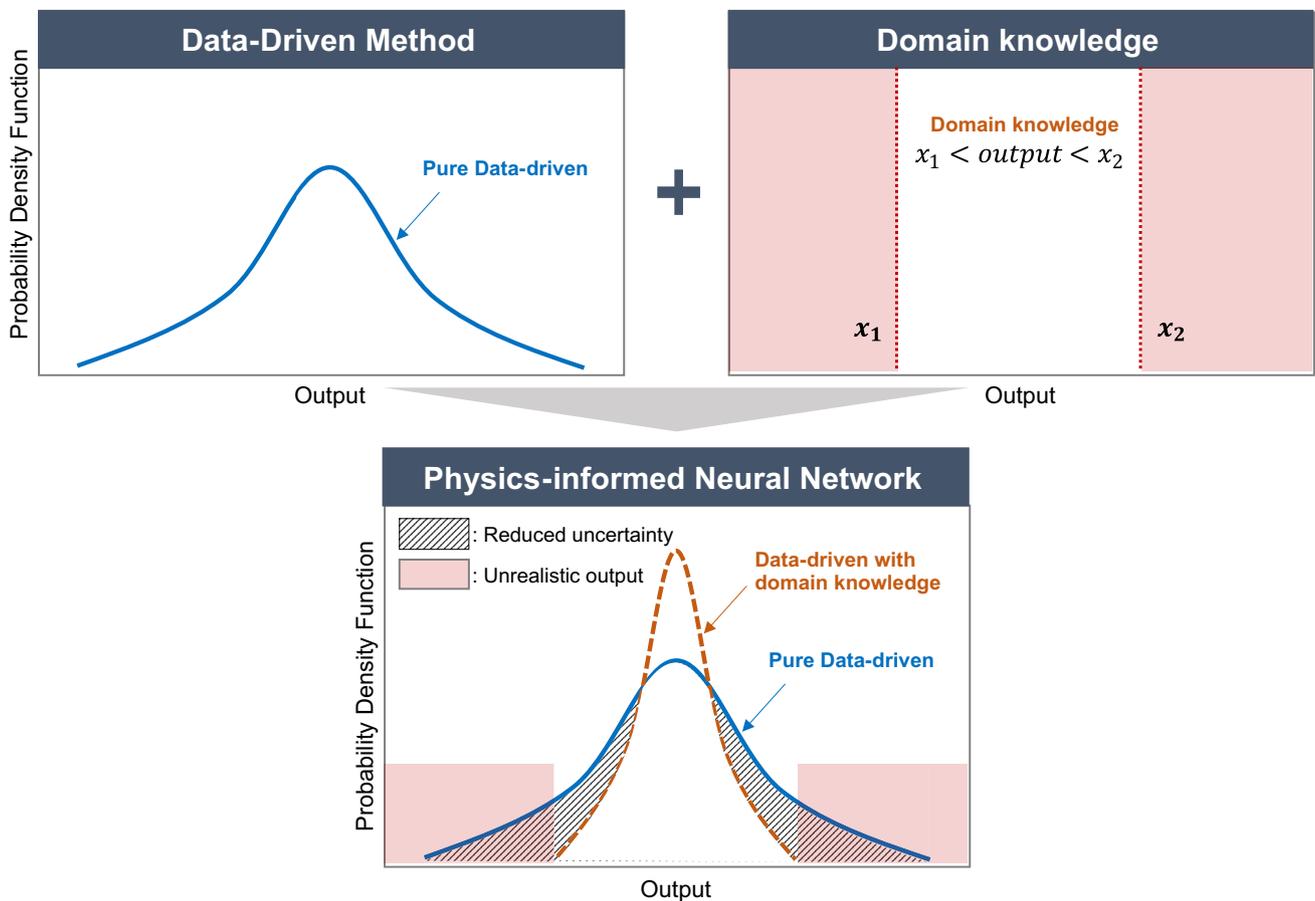


Fig. 1 Uncertainty reduction in the data-driven method using domain knowledge

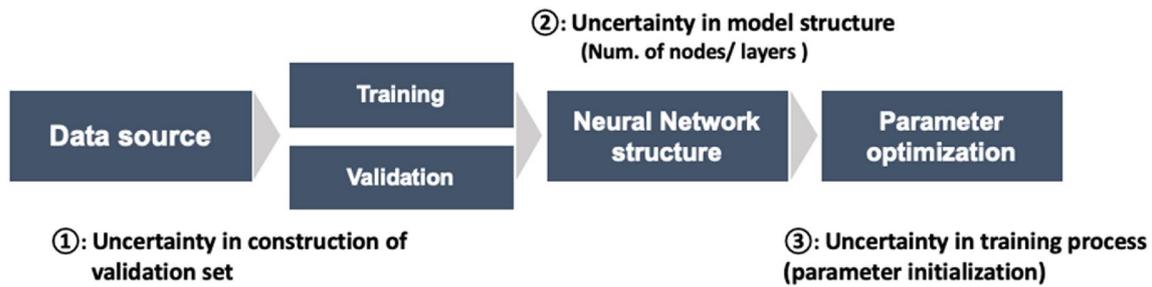


Fig. 2 Training uncertainty in neural network

to avoid overfitting and remove the uncertainty in data selection. Different from the early stopping criterion, it is possible to utilize all data for the training set. Traditionally, the constraint term has been defined as the square sum of model parameters (i.e., weights) (Goodfellow et al. 2016). However, the traditional constraint term does not contain any physical meaning. In this paper, a physics-based constraint term is proposed based on domain knowledge in the extrapolation region. As a result, it is possible to remove the uncertainty in selecting the validation set and make full use of training data.

2.2.2 Data uncertainty

Condition monitoring data that are used for PHM contain noise due to measurement environment or unknown factors. Since the nature of the noise is random, the data will be different every time they are measured. This type of uncertainty in measured data is aleatory uncertainty, and they are irreducible unless the measurement method is changed. In order to simulate the effect of random noise,

often accurate data are generated using a physics model, and then, random noise is added to the data. Figure 3 illustrates that measured data can be different even under the identical true degradation function and the same statistical distribution of noise. In this case, two different sets of random noise, n , are generated from a uniform distribution between ± 0.003 and are added to the same degradation curve (black solid line). Even if the two sets of data are generated from the same distribution of noise, Figs. 3a and b visually show quite different data. In Fig. 3a, measured data seem to be close to the ground truth and the noise looks small between 4000 and 6000 cycles, whereas Fig. 3b has evenly distributed noise. This figure implies that the performance of NN can be different for different sets of training data, which is referred to as data uncertainty in this paper. To quantify the data uncertainty and investigate the robustness of the NN model, in this paper, the performance of the NN is evaluated for different sets of noise (i.e., different random samples of n with the same level of noise).

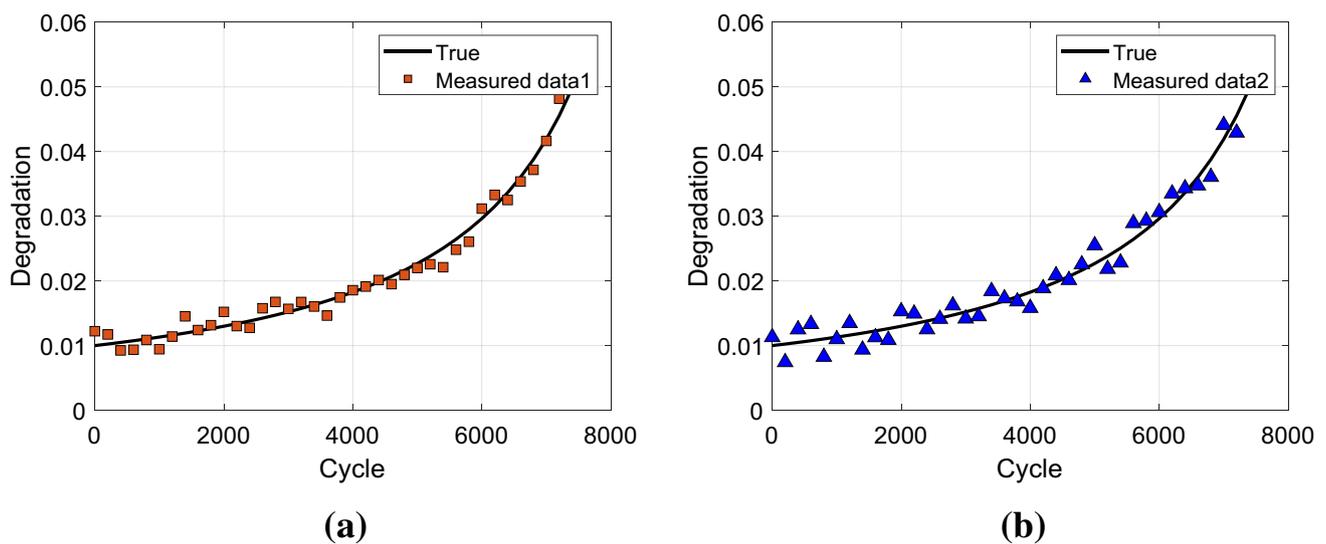


Fig. 3 Data uncertainty in measurement: **a** first noise perturbation **b** second noise perturbation

3 Proposed method

Prognostics aims to predict the future trend based on information from the past. This can be considered as predicting the extrapolation region (future time) using training data in the interpolation region (past time). Since the training process identifies the model parameters of NN (i.e., weights and biases) by minimizing the error between predictions and training data within the interpolation region, prediction uncertainty dramatically increases in the unexplored extrapolation region where training data cannot cover. The primary reason for this uncertainty is the lack of data, which causes NN predictions not to satisfy the underlying physics. To reduce the prediction uncertainty in NN, it is necessary to guide the training process for the extrapolation region with physical knowledge. Therefore, the main idea of the proposed method is to train the NN with both training data in the interpolation region and physical knowledge for the extrapolation region simultaneously. Especially, the physical knowledge used in this method is low-fidelity physical information that refers to level 1 introduced in Sect. 2.1.

Figure 4 shows the basic concept of the proposed approach. Assume that there are training data whose input data and label are \mathbf{x} and \mathbf{y} , respectively. Then, the NN whose weights and biases are θ can calculate the interpolation error, MSE , between the output, $\hat{\mathbf{y}}$, and label, \mathbf{y} . The NN can also predict the outputs in the extrapolation region, $\hat{\mathbf{y}}_p$, whose input is denoted as \mathbf{x}_p based on the current weights and biases. Different from the interpolation

error, MSE , it is impossible to obtain the error for the extrapolation region since there is no available data in this region. Instead, the predictions can be penalized if they violate low-fidelity physical information. This type of penalty function is expressed as $phy(\cdot)$. The cost function for the training process is defined by combining the interpolation error and physical constraint as

$$\text{cost} = MSE + \lambda \cdot phy(\hat{\mathbf{y}}_p), \tag{1}$$

where λ represents the Lagrange multiplier. The objective is to train the NN parameters in the manner of supervised learning with training data as well as to minimize the violation of physical knowledge.

The readers might realize that the physical constraint term may act as a role of a regularization term such as L_1 and L_2 norms of parameters, which reduces the complexity of the NN model due to overfitting. In the same manner, the physical constraint term ensures that the NN is generalized better for the data by preventing overfitting. An advantage of the constraint based regularization for the NN is that it can utilize the whole training data set (Goodfellow et al. 2016). As a result, the proposed method will find the optimum weights and biases, θ , that minimize the cost function while satisfying the physical knowledge. It should be noted that it is possible to use different optimization algorithms to find the optimal θ depending on the form of physical constraint derived from the physical knowledge.

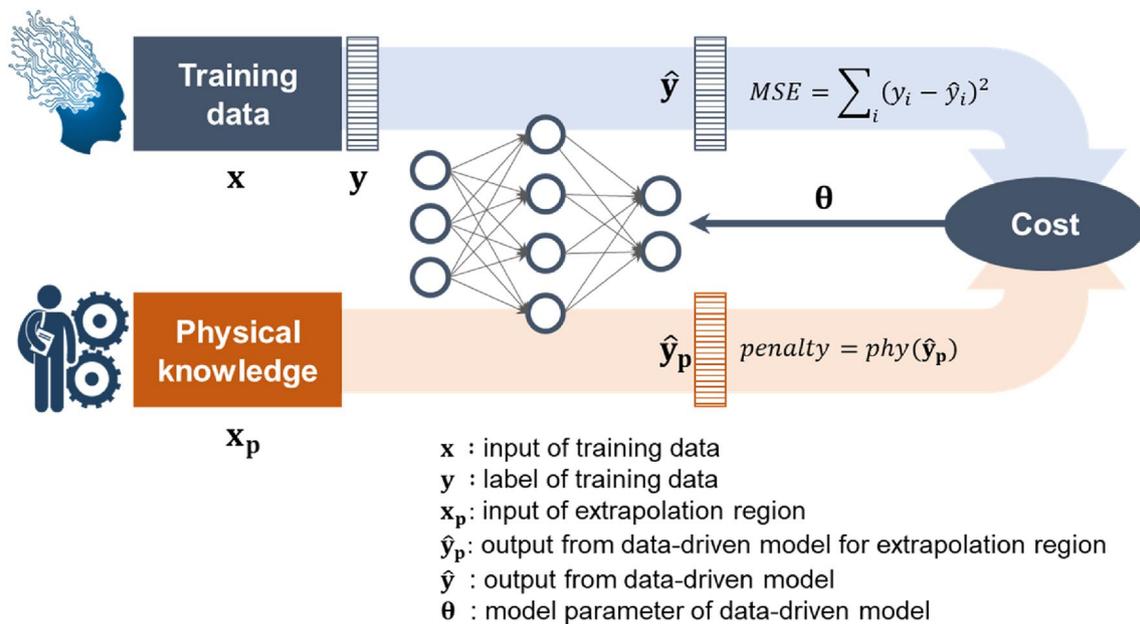


Fig. 4 The framework of physics-informed neural network

4 Case studies

In this paper, we employed two case studies with different purposes. In the first case study, the effect of the proposed method in uncertainty reduction is demonstrated using cooling fan run-to-fail data. Synthetic crack propagation data generated from the Paris law are used for the second application. This study shows the robustness of the proposed method for the uncertainty in test data. The common assumption for both studies was that there is no historical run-to-fail data to train data-driven methods, but only domain knowledge exists. This assumption makes the conventional approaches based on curve fitting or extrapolation since they require the empirical degradation model, which can be obtained from the historical run-to-fail data.

4.1 Cooling fan bearing degradation

Cooling fans are important components for air circulation and temperature control in many industrial fields. Its malfunction or degradation of cooling efficiency can lead to serious damage to the core systems (Peng and Su 2021). In this case study, the proposed method is applied to the prognostics of cooling fan bearing. For this purpose, experiments were conducted to obtain the cooling fan run-to-fail data.

4.1.1 Experimental setup

For the experiment, the DC cooling fan unit was mounted on the test rig with four bolts, and an accelerometer is installed next to the bearing. The specification of the bearing, NSK 693ZZ, is given in Table 2. Figure 5a shows the cooling fan bearing test rig in which the run-to-fail data are obtained. For the accelerated life test, a bolt is mounted on one of the blades to increase the imbalance force, and the cooling fan was operated at an elevated temperature of 70°C in the chamber shown in Fig. 5b. Acceleration signals were obtained every 10 min for 2.5 s, which is counted as one cycle in this paper. The cooling fan was operated at 3000 rpm and the experiment was stopped when the kurtosis of signals reached 4 g since the kurtosis value of normal bearing is equal to 3 g (Kim et al. 2020b). In Fig. 5c, acceleration signals that were obtained for about 120 h are shown. To implement an effective prognostics algorithm, it is important to extract a feature that reflects the

Table 2 Cooling fan bearing specification

Bearing type	Inner diameter (mm)	Outer diameter (mm)	Width (mm)	Dynamic load (N)	Maximum RPM
NSK 693ZZ	3	8	4	560	60,000

health condition of the cooling fan. Among various types of features, we defined a health index of the bearing as the root-mean-square (RMS) of acceleration signals at each cycle, since it is one of the widely used statistical measures that are effective for bearing condition monitoring, and its trend over time is shown in Fig. 5d. It should be noted that the development of an effective health index for condition monitoring is out of the scope of this paper. In the real application, any features or health index can be used for the proposed PINN.

4.1.2 Implementation of PINN

In this case study, we assumed that there is no historical run-to-fail data but domain knowledge about the bearing degradation behavior. In other words, the ANN will be used to predict future behavior without historical run-to-fail data. Two kinds of physical knowledge are employed for this case study. First, the bearing degradation process is divided into three stages: the normal condition, smooth wear condition, and severe wear condition before failure (Shi et al. 2021). This knowledge is adopted by using time t , t^2 and t^3 as inputs for the PINN to predict the health index \hat{y}_t at t . This is because the third-order polynomial function contains two stationary points that can be interpreted as dividers of three stages. Time t means the cycle of cooling fan bearing data. Second, degradation of the bearing cannot be reversed. It means that the health condition of the bearing should follow a monotonic increasing/decreasing trend. Those ANN parameters (weights and biases) that predict the degradation in the extrapolation region (future time) not following this trend should be eliminated or penalized. This physical knowledge is imposed in the following form of a constraint:

$$\sum_{i=2}^{T_p} \max[0, -\Delta\hat{y}_{t_i}] \leq 0, \quad (2)$$

where \hat{y}_{t_i} is the prediction of health index at t_i , and $\Delta\hat{y}_{t_i} = \hat{y}_{t_i} - \hat{y}_{t_{i-1}}$. The physical constraint term is violated whenever the predicted health index \hat{y}_{t_i} becomes smaller than the previous one $\hat{y}_{t_{i-1}}$, which represents monotonicity. It should be noted that the number of prediction points, T_p , should be defined prior to the training process to calculate the physical penalty. Finally, the training of PINN can be interpreted as a constrained optimization problem as shown in the following equation:

$$\begin{aligned} \text{Minimize} \quad & MSE = \sum_{i=1}^N (\hat{y}_i - y_i)^2 \\ \text{Subject to} \quad & \sum_{i=2}^{T_p} \max[0, -\Delta\hat{y}_{t_i}] < 0 \end{aligned} \quad (3)$$

This process can be explained visually in Fig. 6. As shown in the figure, at 250 cycles, N number of training data exist, and they are used to train the NN. During each iteration of the training process, the physical constraint, $\Delta\hat{y}_{t_i}$

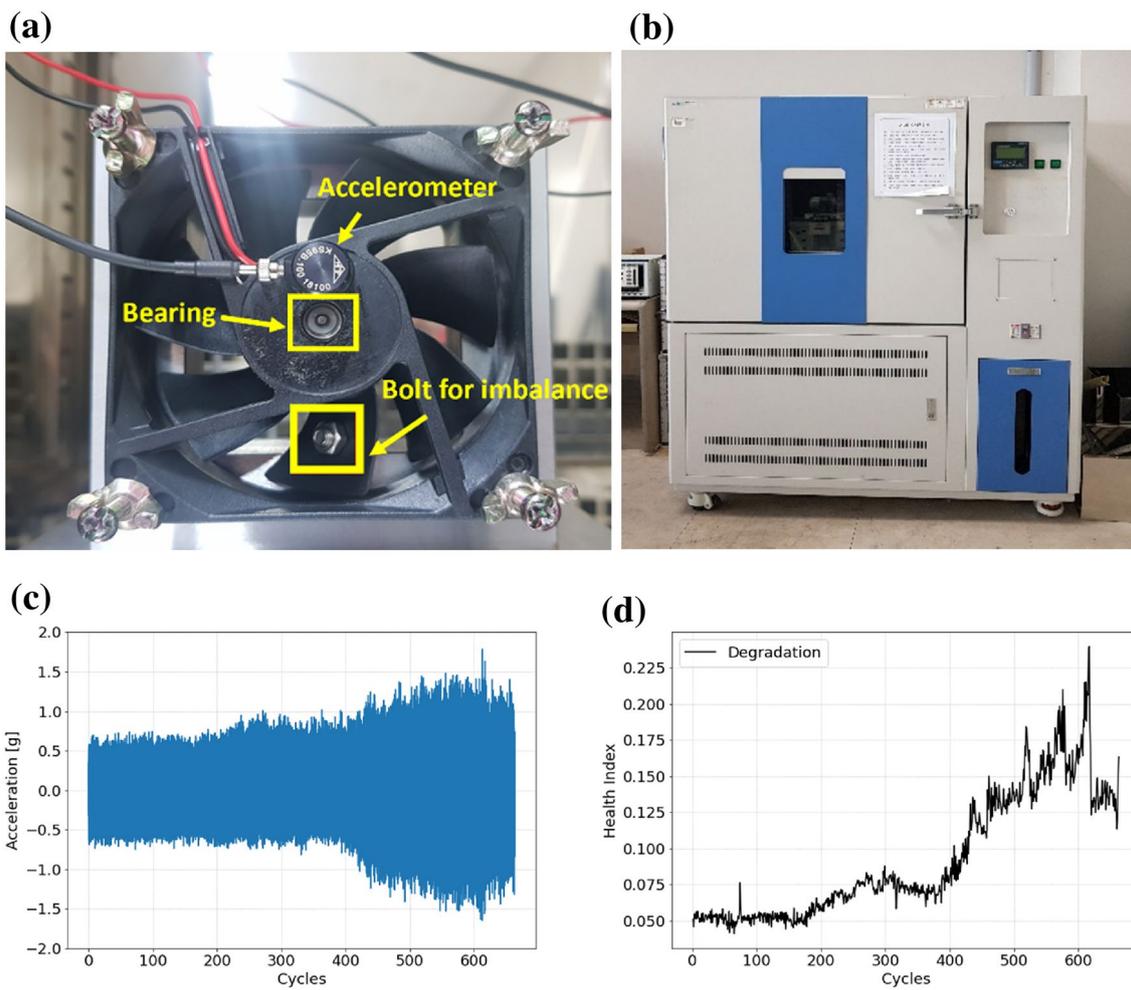


Fig. 5 Experimental apparatus and measurement data of cooling fan bearing: **a** close-up view of the experimental setup, **b** environmental chamber, **c** acceleration signal of the cooling fan, and **d** RMS trend

($i = 1, 2, \dots, T_p$), is calculated, and those training parameters that show decreasing trends will be penalized. For example, the blue dotted line shows the prediction from the NN, and

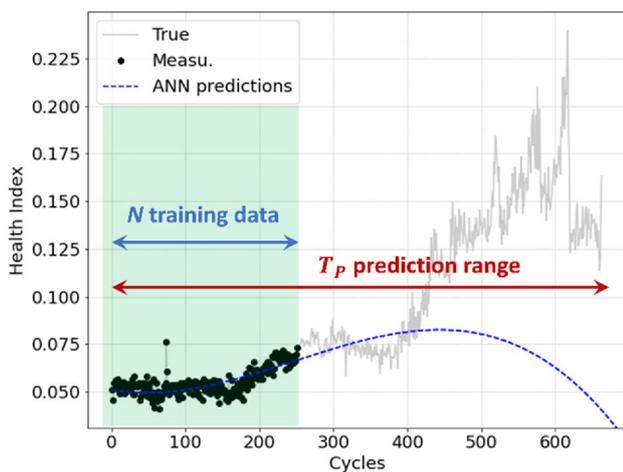


Fig. 6 Physical penalty calculation

its prediction has a decreasing trend after 450 cycles, which means that current weights and biases will be penalized since they violate the physical trend.

To solve the inequality constrained optimization problem, the constrained optimization formulation can be approximated into an unconstrained problem by using a Lagrange multiplier method (Kim et al. 2001). As a result, Eq. (3) can be converted into the following cost function:

$$\cos t = MSE + \lambda_1 \sum_{i=2}^{T_p} \max[0, -\Delta \hat{y}_t], \tag{4}$$

where λ_1 is the Lagrange multiplier for the physical constraint term. Finally, the PINN obtains optimum weights and biases for minimizing the error related to training data, while satisfying the physical information for the extrapolation region. Figure 7 illustrates the framework of PINN for cooling fan bearing prognostics. Since three input nodes are used in the input layer, only two hidden nodes were used to simplify the network structure.

4.1.3 Result of the PINN

To demonstrate the effect of uncertainty reduction of the proposed method, this case study compared three different NNs-based predictions: (1) ordinary NN, (2) PINN with training uncertainty, and (3) PINN without training uncertainty. The main difference between (2) and (3) models is training uncertainty. For the PINN with training uncertainty, existing data were divided into training and validation sets the same as the ordinary NN, and only training data were used for model training. The PINN without training uncertainty represents the proposed method that utilizes all the existing data for model training. To quantify the uncertainty in the training process, 30 NNs with the same structure are trained by three approaches.

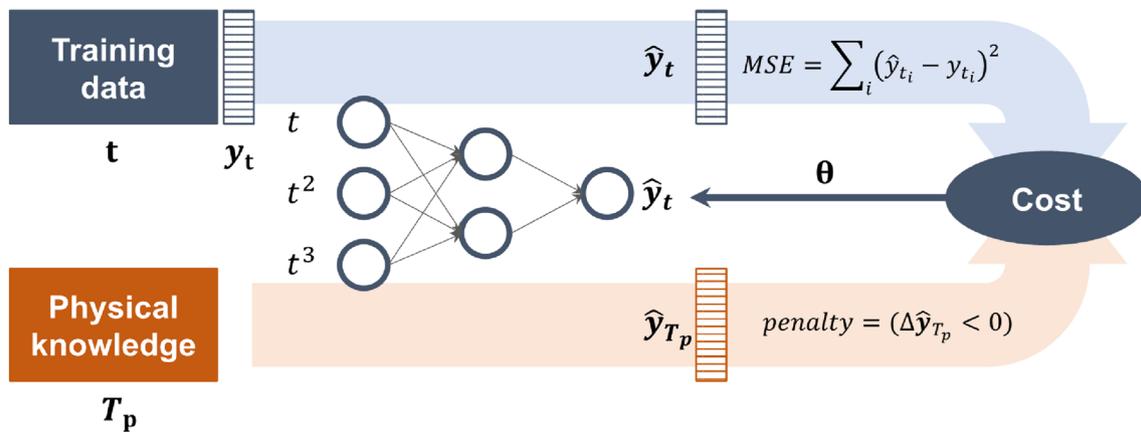
Figure 8 shows the results from these three different NNs. Blue dashed lines in Figs. 8a–c are the predictions from 30 NNs at 200 cycles. Figure 8a indicates that many outputs from the ordinary NN do not follow the physical knowledge that the health index should monotonically increase. On the contrary, PINN brings the predictions that meet the domain knowledge, as shown in Figs. 8b and c. The training uncertainty in Fig. 8a is significantly reduced in Figs. 8b and c. It shows that the exploitation of physical constraint terms has an effect of shaping the training uncertainty by guiding the NN to satisfy the domain knowledge. In addition, uncertainty in Fig. 8c becomes narrower than that of Fig. 8b since all existing data are used for model training. In fact, the uncertainty in Fig. 8c is almost ignorable. Figures 8d, e,

and f illustrate the same results when data up to 480 cycles were used for training. The NNs with physical constraint terms still show better performance in terms of uncertainty than the ordinary NN. This result suggests two important aspects. First, the current prediction uncertainty only represents the confidence of the model, which means that the predictive interval of the ordinary NNs will be larger when the measurement noise is combined. It will hinder NN-based decision-making. Second, even if the uncertainty in the ordinary NNs covers the actual degradation behavior, it should be considered ‘fake’ uncertainty since it includes unrealistic results.

4.2 Crack growth problem

4.2.1 Implementation of PINN

To investigate the robustness of the proposed method, the crack growth example that is widely used in the field of prognostics research is employed. In this approach, synthetic data are first generated from the physics model with the true model parameters, and then, random noise is added to them in order to simulate the actual measurement environment. The information of the physics model and parameters are only used for the purpose of generating data. To generate the synthetic degradation data, Paris’ law (Paris and Erdogan 1960) that describes the crack propagation is employed as follows:



- x : input of training data
- y : label of training data
- x_p : input of extrapolation region
- y_p : output from data-driven model for extrapolation region
- \hat{y} : output from data-driven model
- θ : model parameter of data-driven model

Fig. 7 Physics-informed neural network structure for cooling fan prognostics

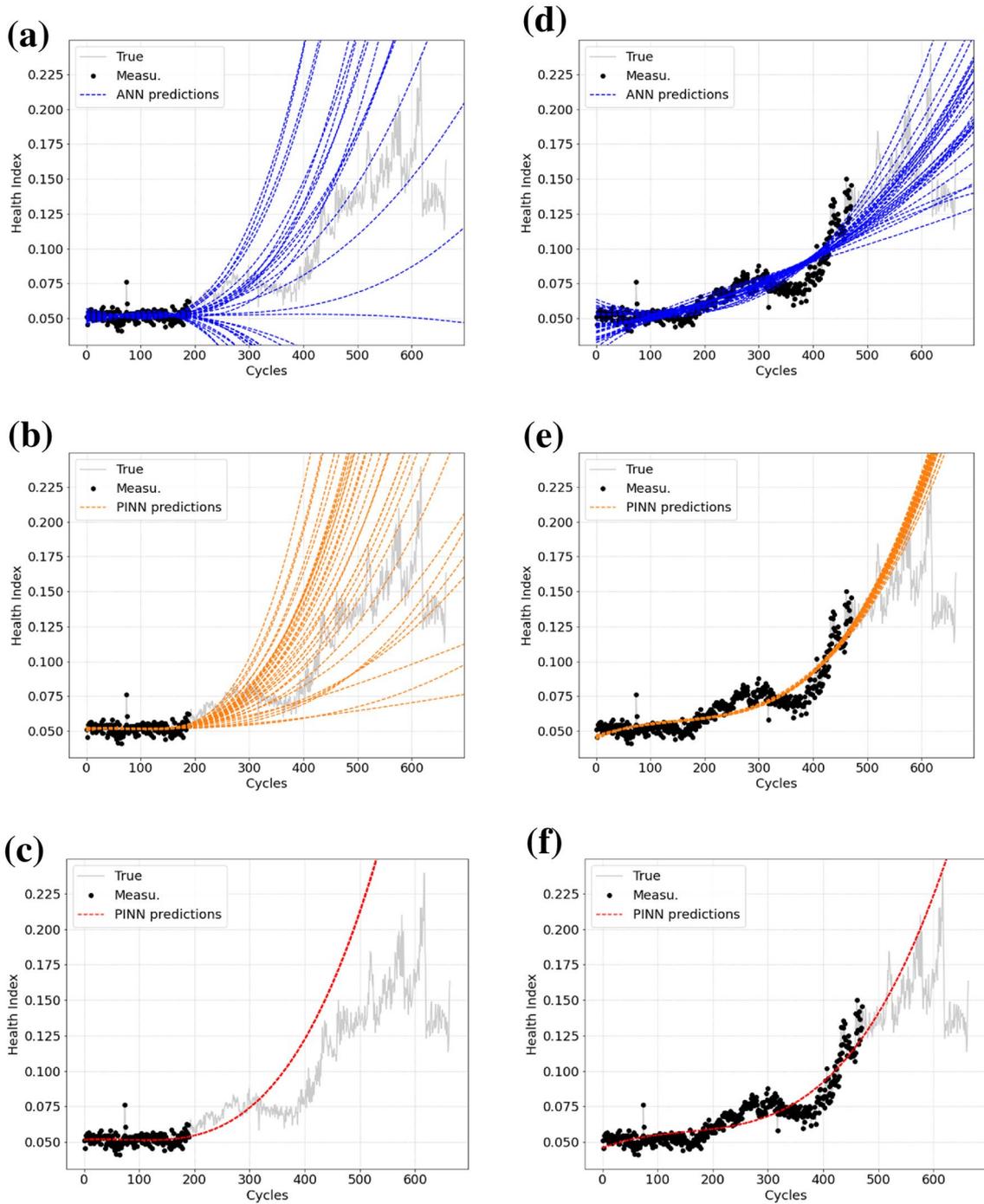


Fig. 8 Prediction results: **a** ordinary NN at 200 cycles, **b** PINN with training uncertainty at 200 cycles, **c** PINN without training uncertainty at 200 cycles, **d** ordinary NN at 480 cycles, **e** PINN with training uncertainty at 480 cycles, **f** PINN without training uncertainty at 480 cycles

$$\frac{da}{dN} = C(\Delta K)^m, \Delta K = \Delta\sigma\sqrt{\pi a} \tag{5}$$

$$z_t = a_t + n, \tag{6}$$

where N is the number of cycles, a is the half crack size, z_t is measured crack size, ΔK is the range of stress intensity factor, $\Delta\sigma$ is the stress range, and C and m are model parameters. In order to simulate the measurement noise, n , random noise uniformly distributed between $\pm umm$ are introduced ($u = 1$ mm). Figure 9 shows the simulated crack growth data when $C = 1.5 \times 10^{-10}$, $m = 3.8$, $\Delta\sigma = 60MPa$, and the initial crack size, $a_0 = 0.01m$. Assumptions made for this study can be summarized as follows:

- There is no run-to-fail data.
- Measured crack size data are only available (z_t).
- The physical degradation model is unavailable.
- Only Low-fidelity physical information is available.

In summary, the goal is to predict future behavior with low-fidelity physical information without historical data and physical models. It should be noted that the physical model, i.e., Paris' law, is only employed to generate the simulation crack growth dataset. Only the training data in the interpolation region and the low-fidelity physical information are utilized for NN training. Two pieces of low-fidelity physical information are employed: (1) the crack size should not be smaller than the previous one (i.e., monotonic increasing trend), and (2) the crack size will exponentially increase in the future (i.e., non-linear behavior). The prediction range where the low-fidelity physical information is embedded is shown in Fig. 9. Blue and red dotted lines in Fig. 9 show the outputs from the ordinary NN for interpolation and extrapolation region, respectively. Although the ordinary NN is

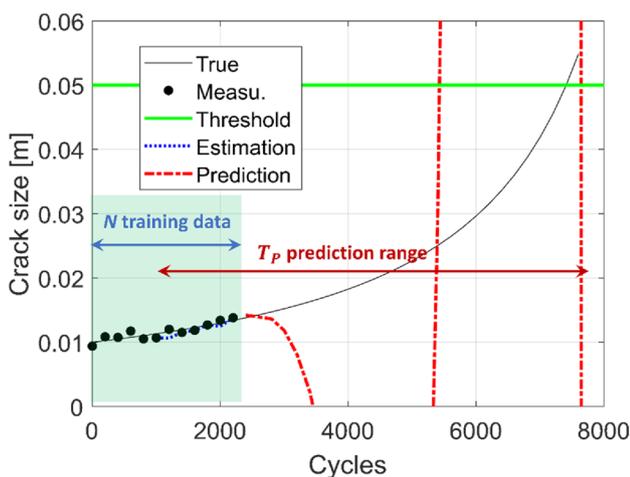


Fig. 9 Physical penalty calculation

trained well within the interpolation region, it shows severe fluctuating predictions that violate the physical knowledge in the extrapolation region. In physics, it is not reasonable that the crack size becomes smaller than 0 and decreases from the previous step. To prevent the NN from yielding physically unreasonable predictions, it is vital to guide the NN in the training process to follow two low-fidelity physical information. For this purpose, the proposed method can be defined as a constrained optimization problem, and its optimization formulation is given as follows:

$$\begin{aligned} \text{Minimize } & MSE = \sum_{i=1}^{N_t} (\hat{y}_i - y_i)^2 \\ \text{Subject to } & \sum_{i=2}^{T_p} \max[0, -\Delta\hat{y}_i] \leq 0, \\ & \sum_{i=2}^{T_p} \max[0, -\Delta(\Delta\hat{y}_i)] \leq 0 \end{aligned} \tag{7}$$

where N_t is the number of training data, $\Delta\hat{y}_i = \hat{y}_i - \hat{y}_{i-1}$, and $\Delta(\Delta\hat{y}_i) = \Delta\hat{y}_i - \Delta\hat{y}_{i-1}$. The first constraint term expresses the monotonic increase of crack size. This physical constraint term is violated whenever the predicted crack size is smaller than the one-step before crack size. The constraint related to the non-linear increasing trend is encoded as the second derivative of crack size. When the slope of crack growth shows a negative trend (i.e., a non-linearly decreasing trend), the term penalizes the current model parameters. It should be noted that the number of prediction points, T_p , should be defined prior to the training step to calculate the physical constraint. The prediction points can include both the interpolation and extrapolation regions. As same in Sect. 4.1.2, the optimization formulation can be converted into an unconstrained problem as follows:

$$\cos t = MSE + \lambda_1 \sum_{i=2}^{T_p} \max[0, -\Delta\hat{y}_i] + \lambda_2 \sum_{i=2}^{T_p} \max[0, -\Delta(\Delta\hat{y}_i)], \tag{8}$$

where λ_1 and λ_2 are Lagrange multipliers for two physical constraint terms. Finally, the PINN obtains optimum weights and biases for minimizing the error related to training data, while satisfying the physical information in the extrapolation region.

4.2.2 Training uncertainty quantification

In this study, the proposed PINN prognostics model is constructed based on a feedforward network with five input nodes and one hidden layer with three nodes. To minimize the complexity of the model, a pure-linear function is used as an activation function for the two layers. The network is trained to predict a one-step-ahead crack size based on the previous five crack sizes, i.e., the five most recent crack sizes, $z_{t-4}, z_{t-3}, \dots, z_t$ are used to predict the \hat{z}_{t+1} when the current time step is represented as t . To quantify the training uncertainty in the NN, 50 NNs with the same structures are

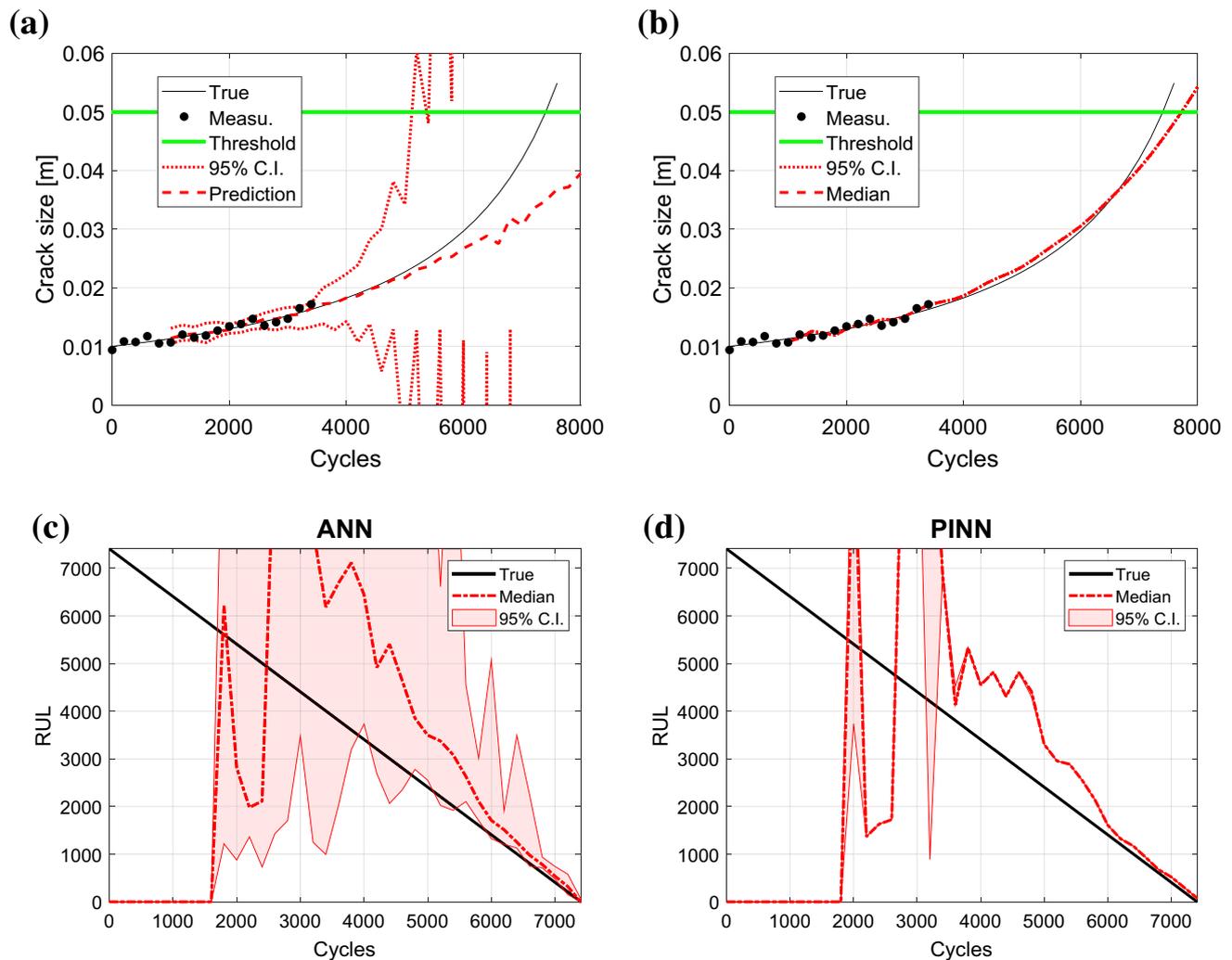


Fig. 10 Prediction of crack growth: **a** ordinary neural network, **b** proposed method, **c** RUL estimation of the ordinary neural network, and **d** RUL estimation of the proposed method

trained. Figures 10a and b are the prediction results obtained from the ordinary NN (i.e., artificial neural network; ANN) and the proposed method (PINN), respectively. Dashed and dotted lines indicate the median and 95% confidence interval (C.I.) of predictions out of 50 ANN predictions. The prediction result from the ANN contains large uncertainty and a highly fluctuating trend that violates the physics in the extrapolation region. When the physical knowledge is embedded into the NN, prediction uncertainty is dramatically reduced, and the median of prediction becomes closer to the ground truth compared to the ANN. This phenomenon is notable in RUL prediction as shown in Figs. 10c

and d. The median of RUL prediction from PINN has narrow uncertainty and is close to the ground truth of RUL. Although there is a bias between prediction and true RUL, the proposed method shows improved performance in terms of uncertainty considering the initial assumption that there was no historical run-to-fail data.

4.2.3 Data uncertainty quantification

As mentioned in Sect. 2, it is crucial for the data-driven method to have consistent performance for the different test datasets. To investigate the robustness of the proposed

method, 10 different noise perturbations drawn from the same uniform distribution are added to one crack growth trend. In other words, 10 sets of measurements with the same noise level are generated. Then, the proposed method and the ANN are applied to the individual run-to-fail dataset to predict the RUL. Per each dataset, 50 network models are trained for training uncertainty quantification, and lower (2.5 percentile), median (50 percentile), and upper bound (97.5 percentile) of RUL predictions are stored. Finally, it is possible to obtain 10 sets of 95% C.I. and the median of the RUL predictions, as shown in Table 3. Uncertainty in a row of the table indicates the epistemic uncertainty that can be reduced by gathering more data or improving model quality. In contrast, the column corresponds to the aleatory uncertainty derived from the data's

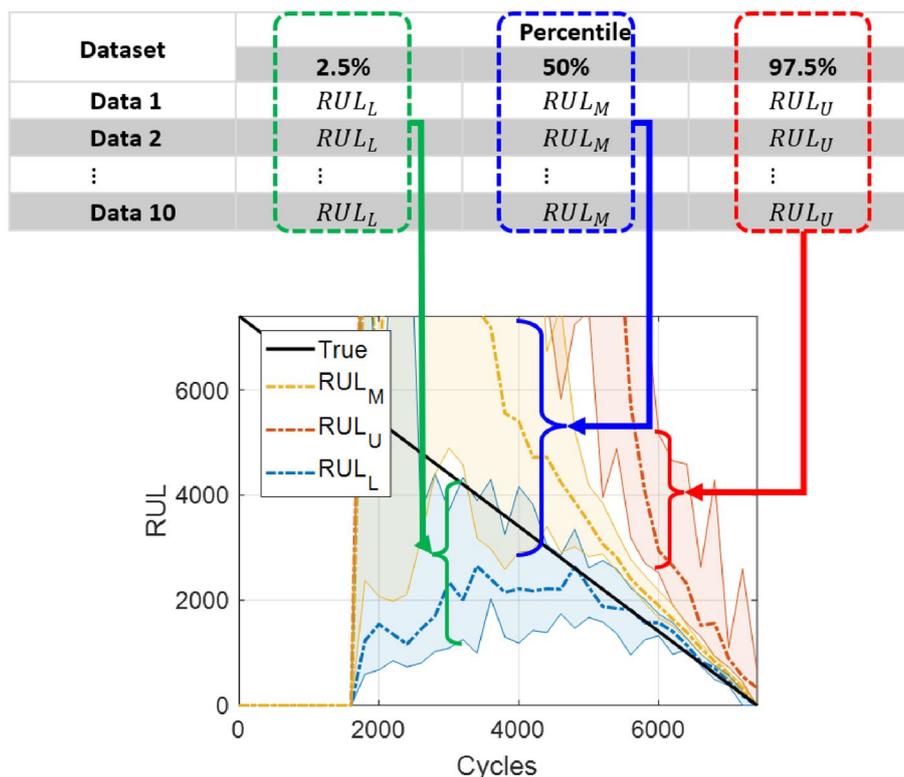
randomness. To represent the effect of data uncertainty on the RUL prediction, 95% C.I. and median for 10 sets of RUL_L , RUL_M , and RUL_U are calculated. To help the understanding of readers, the procedure of data uncertainty quantification is shown in Fig. 11. In the figure, each color represents the 95% C.I. and median of RUL prediction. The shaded surface is the uncertainty bounds for RUL estimation due to data uncertainty at each cycle. Figures 12a and b are data uncertainty in RUL prediction obtained from the ANN and the proposed method, PINN, respectively. Figure 12a shows that the uncertainty in the ANN prediction dramatically increases when the data uncertainty is considered. In the real application, the overall uncertainty range should be defined as the distance from the lower bound of RUL_L and the upper bound of RUL_U . It means the appropriate decision-making cannot be established based on the ANN prognostics model. On the other hand, overall uncertainty in RUL prediction of the PINN is reduced to a narrower bound, as shown in Fig. 12b. Moreover, medians of RUL_L , RUL_M , and RUL_U are close to each other, which means that the RUL prediction is more robust against the data uncertainty (i.e., measurement noise).

To identify the robustness of the proposed method, two different levels of noise ($u = 2mm$ and $3mm$) are also

Table 3 RUL prediction for 10 different noise perturbations

Dataset	Percentile		
	2.5%	50%	97.5%
Data 1	RUL_L	RUL_M	RUL_U
Data 2	RUL_L	RUL_M	RUL_U
⋮	⋮	⋮	⋮
Data 10	RUL_L	RUL_M	RUL_U

Fig. 11 Data uncertainty quantification process



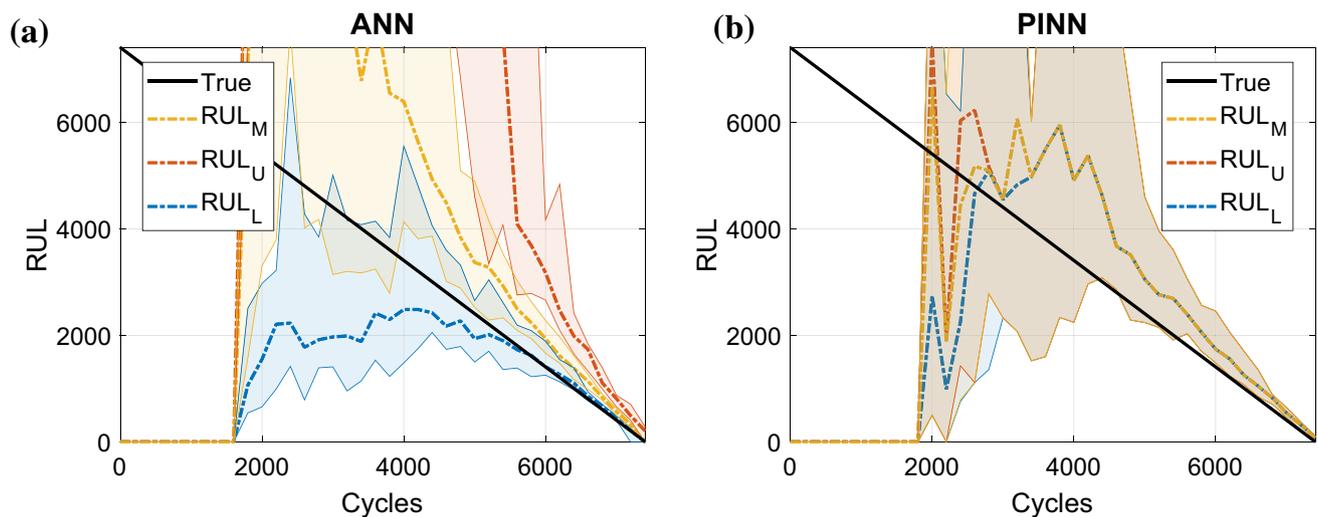


Fig. 12 Data uncertainty quantification in RUL prediction: **a** ordinary neural network, and **b** the proposed method

considered for the crack growth problem. Figure 13 shows the comparison study between the ANN and the PINN for different levels of noise. In general, the RUL prediction uncertainty increases as a higher noise level is added to the true crack growth behavior. For the ANN models, however, the uncertainty in RUL prediction does not converge even until 6000 cycles, which corresponds to about 85% of end of life. In contrast, the uncertainty range of PINN was about 2000 cycles at 6000 cycles for both cases. It shows that the proposed method is more robust than the ANN-based prognostics in the absence of run-to-fail data.

5 Conclusions and future works

In applying prognostics to the real industry, there are several practical challenges related to the lack of a high-fidelity physical model and run-to-fail data. The utilization of advantages of both physics-based and data-driven prognostics methods is demanding to address these problems. This paper proposed a PINN-based prognostics method that trains and guides the NN in the interpolation and extrapolation

region at the same time. As a result, the model parameters of the NN are optimized to minimize the interpolation error, while satisfying the physical knowledge in the extrapolation region. The results from two case studies show that the proposed method shapes the uncertainty in data-driven predictions and brings more robust predictions than ordinary NN. Although the proposed method shows an improved result in terms of uncertainty, there are still limitations that should be explored in the future. The RUL predictions from the PINN are biased from the ground truth since the case studies in this paper assumed that there is no available run-to-fail data. This also happened from the physics that the degradation rate is accelerated toward the end of life. This bias should be modeled and considered in prognostics to achieve effective decision-making. Future work will be focused on estimating the bias and improving the prognostics accuracy in real time. Moreover, it would be beneficial to develop the method to integrate the physical knowledge with various types of data-driven prognostics algorithms such as the gaussian process, support vector machine, or regression. In future work, various types of physics-informed data-driven prognostics will be developed, and a comparison study will be conducted.

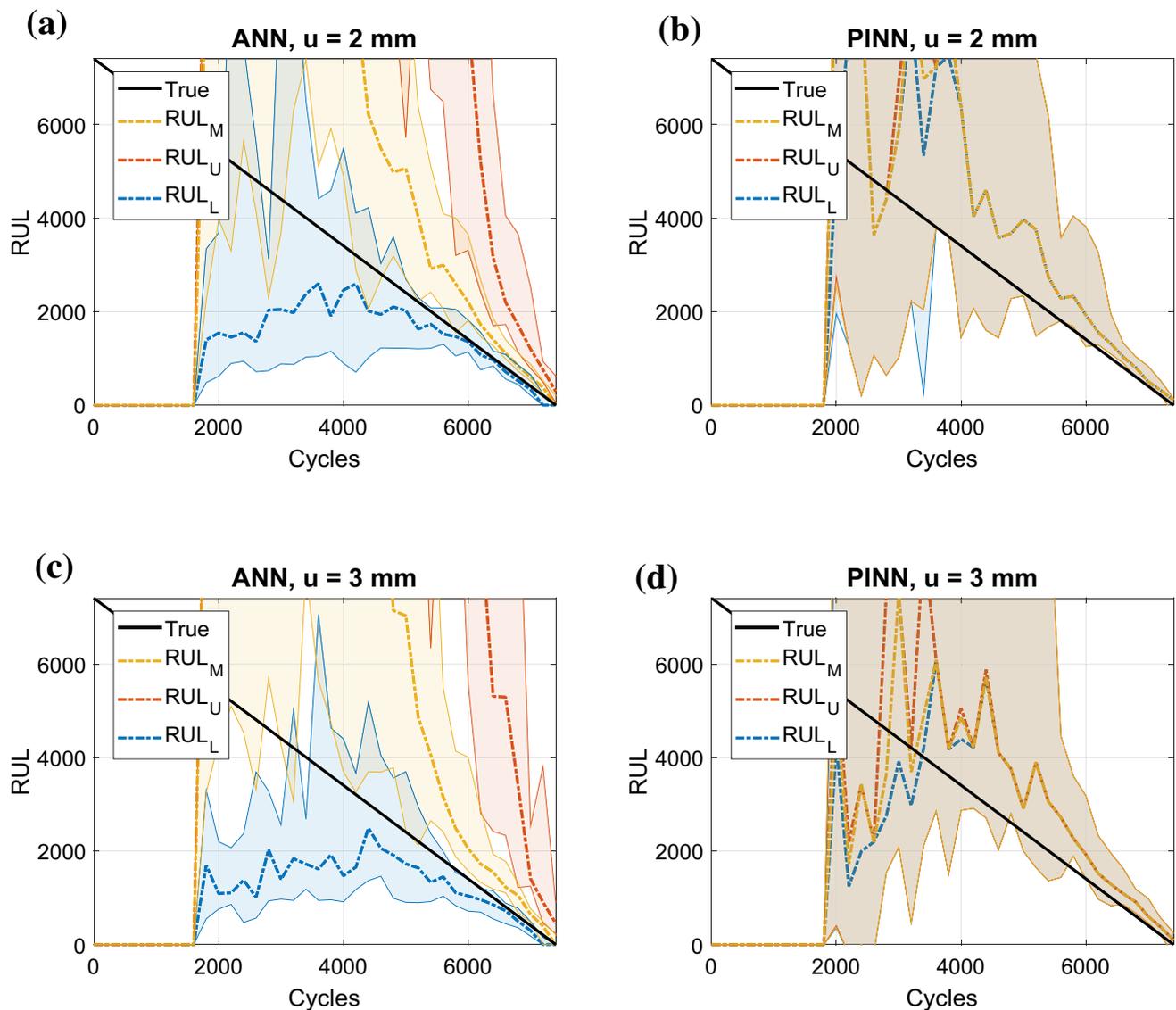


Fig. 13 Comparison of prognostics performance between different levels of noise: **a** RUL prediction of ANN ($u=2$ mm), **b** RUL prediction of PINN ($u=2$ mm), **c** RUL prediction of ANN ($u=3$ mm), and **d** RUL prediction of PINN ($u=3$ mm)

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1007/s00158-022-03348-0>.

Acknowledgements This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korean government (MSIT) (No. 2020R1A4A4079904).

Funding National research foundation of Korea, No. 2020R1A4A4079904, Joo-Ho Choi.

Declarations

Conflict of interest The authors have no competing interests to declare that are relevant to the content of this article.

Replication of Results The python script of the cooling fan in Sect. 4.1 is included in the supplementary document along with data file bearing_hi.csv.

References

- Dourado A, Viana FAC (2020) Physics-informed neural networks for missing physics estimation in cumulative damage models: a case study in corrosion fatigue. *J Comput Inf Sci Eng* 20(6):61007
- Fang Z, Zhan J (2019) A physics-informed neural network framework for PDEs on 3D surfaces: time independent problems. *IEEE Access* 8:26328–26335

- Goodfellow I, Bengio Y, Courville A (2016) *Deep learning*. MIT press
- Goswami S, Anitescu C, Chakraborty S, Rabczuk T (2020) Transfer learning enhanced physics informed neural network for phase-field modeling of fracture. *Theor Appl Fract Mech* 106:102447
- Ham S, Han SY, Kim S, Park HJ, Park KJ, Choi JH (2019) A comparative study of fault diagnosis for train door system: traditional versus deep learning approaches. *Sensors* 19(23):5160. <https://doi.org/10.3390/s19235160>
- Heimes FO (2008) Recurrent neural networks for remaining useful life estimation. *Progn Heal Manag PHM 2008 Int Conf*. <https://doi.org/10.1109/PHM.2008.4711422>
- Huang X, Torgeir M, Cui W (2008) An engineering model of fatigue crack growth under variable amplitude loading. *Int J Fatigue*. <https://doi.org/10.1016/j.ijfatigue.2007.03.004>
- Kim S, Choi JH (2019) Convolutional neural network for gear fault diagnosis based on signal segmentation approach. *Struct Heal Monit* 18(5–6):1401–1415. <https://doi.org/10.1177/1475921718805683>
- Kim NH, Choi KK, Chen JS (2001) Die shape design optimization of sheet metal stamping process using meshfree method. *Int J Numer Methods Eng* 51(12):1385–1405
- Kim S, Kim NH, Choi JH (2020a) Information value-based fault diagnosis of train door system under multiple operating conditions. *Sensors (switzerland)*. <https://doi.org/10.3390/s20143952>
- Kim S, An D, Choi J-H (2020b) Diagnostics 101: a tutorial for fault diagnostics of rolling element bearing using envelope analysis in MATLAB. *Appl Sci* 10(20):7302
- Kim S, Choi J-H, Kim NH (2021) Challenges and opportunities of system-level prognostics. *Sensors* 21(22):7655. <https://doi.org/10.3390/s21227655>
- Lee J, Wu F, Zhao W, Ghaffari M, Liao L, Siegel D (2014) Prognostics and health management design for rotary machinery systems—Reviews, methodology and applications. *Mech Syst Signal Process* 42(1–2):314–334
- Lei Y, Li N, Guo L, Li N, Yan T, Lin J (2018) Machinery health prognostics: a systematic review from data acquisition to RUL prediction. *Mech Syst Signal Process*. <https://doi.org/10.1016/j.ymssp.2017.11.016>
- Lim KYH, Zheng P, Chen C-H (2020) A state-of-the-art survey of Digital Twin: techniques, engineering product lifecycle management and business innovation perspectives. *J Intell Manuf* 31(6):1313–1337
- Mao Z, Jagtap AD, Karniadakis GE (2020) Physics-informed neural networks for high-speed flows. *Comput Methods Appl Mech Eng* 360:112789
- Meng X, Li Z, Zhang D, Karniadakis GE (2020) PPINN: Parareal physics-informed neural network for time-dependent PDEs. *Comput Methods Appl Mech Eng* 370:113250
- Nascimento RG, Viana FAC (2019) Fleet prognosis with physics-informed recurrent neural networks. arXiv Preprint arXiv. <https://doi.org/10.12783/shm2019/32301>
- Negri E, Pandhare V, Cattaneo L, Singh J, Macchi M, Lee J (2021) Field-synchronized digital twin framework for production scheduling with uncertainty. *J Intell Manuf* 32(4):1207–1228
- Paris PC, Erdogan F (1960) A critical analysis of crack propagation laws. *J Basic Eng* 85:528–534
- Peng C-C, Su C-Y (2021) Modeling and parameter identification of a cooling fan for online monitoring. *IEEE Trans Instrum Meas* 70:1–14
- Giorgiani do Nascimento R, Viana F, Corbetta M, Kulkarni CS “Usage-based Lifing of Lithium-Ion Battery with Hybrid Physics-Informed Neural Networks,” *AIAA AVIATION 2021 FORUM*, 2021, p. 3046
- Raissi M, Perdikaris P, Karniadakis GE (2019) Physics-informed neural networks: a deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *J Comput Phys* 378:686–707
- Shi J, Yu T, Goebel K, Wu D (2021) Remaining useful life prediction of bearings using ensemble learning: the impact of diversity in base learners and features. *J Comput Inf Sci Eng*. <https://doi.org/10.1115/1.4048215>
- Wang T (2010) Trajectory similarity based prediction for remaining useful life estimation *Network*
- Yang L, Zhang D, Karniadakis GE (2020) Physics-informed generative adversarial networks for stochastic differential equations. *SIAM J Sci Comput* 42(1):A292–A317
- Yucesan YA, Viana FAC (2020) ‘A physics-informed neural network for wind turbine main bearing fatigue. *Int J Progn Heal Manag* 11(1):17
- Yucesan YA, Viana FAC (2021) Hybrid physics-informed neural networks for main bearing fatigue prognosis with visual grease inspection. *Comput Ind* 125:103386
- Yucesan YA, Viana FAC (2022) A hybrid physics-informed neural network for main bearing fatigue prognosis under grease quality variation. *Mech Syst Signal Process* 171:108875

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.