



Kent Academic Repository

Wu, Shaomin (2021) *Relevance Vector Regression for Remaining Useful Life Prediction*. In: *Mathematical Modeling in Physical Sciences, Social Sciences and Technology (icmm-21)*, 17-18 December, 2021, Organized by Department of Mathematics, Chaudhary Bansi Lal University, 17-18 Dec 2021, Bhiwani, India. (Unpublished)

Downloaded from

<https://kar.kent.ac.uk/92364/> The University of Kent's Academic Repository KAR

The version of record is available from

This document version

Presentation

DOI for this version

Licence for this version

CC0 (Public Domain)

Additional information

Versions of research works

Versions of Record

If this version is the version of record, it is the same as the published version available on the publisher's web site. Cite as the published version.

Author Accepted Manuscripts

If this document is identified as the Author Accepted Manuscript it is the version after peer review but before type setting, copy editing or publisher branding. Cite as Surname, Initial. (Year) 'Title of article'. To be published in *Title of Journal*, Volume and issue numbers [peer-reviewed accepted version]. Available at: DOI or URL (Accessed: date).

Enquiries

If you have questions about this document contact ResearchSupport@kent.ac.uk. Please include the URL of the record in KAR. If you believe that your, or a third party's rights have been compromised through this document please see our [Take Down policy](https://www.kent.ac.uk/guides/kar-the-kent-academic-repository#policies) (available from <https://www.kent.ac.uk/guides/kar-the-kent-academic-repository#policies>).

Relevance Vector Regression for Remaining Useful Life Prediction

Prof. Shaomin Wu

Kent Business School
University of Kent

December 17, 2021

- This presentation is based on the paper: Wang, X., Jiang, B., Wu, S., Lu, N. and Ding, S., 2021. Multivariate Relevance Vector Regression based Degradation Modeling and Remaining Useful Life Prediction. *IEEE Transactions on Industrial Electronics*, doi: 10.1109/TIE.2021.3114724.
- I would like to thank my co-authors of the paper;
- I would also like to thank the organisers of *the International Conference in Mathematical Modeling in Physical Sciences, Social Sciences and Technology (icmm-21)* for their invitations.

Table of Contents

- A brief introduction to Remaining Useful Life and Relevance Vector Regression
- Multivariate Relevance Vector Regression
- A Case Study

Table of Contents

- 1 A brief introduction to Remaining Useful Life and Relevance Vector Regression
- 2 Multivariate Relevance Vector Regression
- 3 Case study

Importance of RUL Prediction and our methods

- **Remaining Useful Life (RUL)** RUL is the length of an industrial item from its current time to the end of its useful life.
- **Importance of RUL estimation** It is needed in condition based maintenance, prognostics and health management.
- **Some methods for RUL estimation**
 - ▶ With lifetime data and other data, one can build a Cox regression model and then derive the probability distribution of the RUL,
 - ▶ One may also estimate the distribution of the first hitting time based on a degradation path, which will be adopted in this talk.

Importance of RUL Prediction and our methods

- **Remaining Useful Life (RUL)** RUL is the length of an industrial item from its current time to the end of its useful life.
- **Importance of RUL estimation** It is needed in condition based maintenance, prognostics and health management.
- **Some methods for RUL estimation**
 - ▶ With lifetime data and other data, one can build a Cox regression model and then derive the probability distribution of the RUL,
 - ▶ One may also estimate the distribution of the first hitting time based on a degradation path, which will be adopted in this talk.

Importance of RUL Prediction and our methods

- **Remaining Useful Life (RUL)** RUL is the length of an industrial item from its current time to the end of its useful life.
- **Importance of RUL estimation** It is needed in condition based maintenance, prognostics and health management.
- **Some methods for RUL estimation**
 - ▶ With lifetime data and other data, one can build a Cox regression model and then derive the probability distribution of the RUL,
 - ▶ One may also estimate the distribution of the first hitting time based on a degradation path, which will be adopted in this talk.

Importance of RUL Prediction and our methods

- **Remaining Useful Life (RUL)** RUL is the length of an industrial item from its current time to the end of its useful life.
- **Importance of RUL estimation** It is needed in condition based maintenance, prognostics and health management.
- **Some methods for RUL estimation**
 - ▶ With lifetime data and other data, one can build a Cox regression model and then derive the probability distribution of the RUL,
 - ▶ One may also estimate the distribution of the first hitting time based on a degradation path, which will be adopted in this talk.

Regression with OLS, Ridge, LASSO, ElasticNet

- **Ordinary least squares (OLS):** Given a set of data points, $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, such that $\mathbf{x}_i (= (1, x_{i1}, \dots, x_{ip})^T) \in R^{n+1}$ is a feature vector for the i th case and $y_i \in R^1$ is a target output, one can build a linear regression model

$$y_n = \mathbf{w}^T \mathbf{x} + \epsilon_n,$$

with the aim to minimise the error on the training datasets

$$\min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \mathbf{w}) \quad (1)$$

where $\mathbf{w} = (w_0, w_1, \dots, w_p)^T$ is the weight vector.

- **Extensions:** Ridge, LASSO (least absolute shrinkage and selection operator), and ElasticNet are extensions of OLS, with an additional **penalty on the weights** that aims to maximise the generalisation.

Regression with OLS, Ridge, LASSO, ElasticNet

- **Ordinary least squares (OLS):** Given a set of data points, $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, such that $\mathbf{x}_i (= (1, x_{i1}, \dots, x_{ip})^T) \in R^{n+1}$ is a feature vector for the i th case and $y_i \in R^1$ is a target output, one can build a linear regression model

$$y_n = \mathbf{w}^T \mathbf{x} + \epsilon_n,$$

with the aim to minimise the error on the training datasets

$$\min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \mathbf{w}) \quad (1)$$

where $\mathbf{w} = (w_0, w_1, \dots, w_p)^T$ is the weight vector.

- **Extensions:** Ridge, LASSO (least absolute shrinkage and selection operator), and ElasticNet are extensions of OLS, with an additional **penalty on the weights** that aims to maximise the generalisation.

Support Vector Regression (SVR)

- **SVR** aims to minimize the weights,

$$\min \frac{1}{2} \mathbf{w}^T \mathbf{w} \quad (2)$$

$$|y_i - \mathbf{x}_i^T \mathbf{w}| \leq \epsilon \quad (3)$$

where \mathbf{w} is the weight vector and ϵ is the error term.

- ϵ -**SVR** is given by

$$\min_{w, b, \xi, \xi^*} \frac{1}{2} w^T w + C \left(\sum_{i=1}^l \xi_i + \sum_{i=1}^l \xi_i^* \right)$$

$$y_i - w^T \phi(x_i) - b \leq \epsilon + \xi_i,$$

$$w^T \phi(x_i) + b - y_i \leq \epsilon + \xi_i^*,$$

$$\xi_i, \xi_i^* \geq 0, i = 1, \dots, l.$$

where C is the regularization parameter and balances the trade-off between the model complexity and empirical error, ξ_i and ξ_i^* are slack variables.

Support Vector Regression (SVR)

- **SVR** aims to minimize the weights,

$$\min \frac{1}{2} \mathbf{w}^T \mathbf{w} \quad (2)$$

$$|y_i - \mathbf{x}_i^T \mathbf{w}| \leq \epsilon \quad (3)$$

where \mathbf{w} is the weight vector and ϵ is the error term.

- ϵ -**SVR** is given by

$$\min_{\mathbf{w}, b, \xi, \xi^*} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \left(\sum_{i=1}^l \xi_i + \sum_{i=1}^l \xi_i^* \right)$$

$$y_i - \mathbf{w}^T \phi(x_i) - b \leq \epsilon + \xi_i,$$

$$\mathbf{w}^T \phi(x_i) + b - y_i \leq \epsilon + \xi_i^*,$$

$$\xi_i, \xi_i^* \geq 0, i = 1, \dots, l.$$

where C is the regularization parameter and balances the trade-off between the model complexity and empirical error, ξ_i and ξ_i^* are slack variables.

Limitations of SVR

- SVR has been used in many applications Nevertheless, they suffer some limitations:
 - ▶ **non-probabilistic:** SVR does not output probabilistic predictions;
 - ▶ **C and ϵ :** parameters of C and ϵ must be determined by cross-validation; and
 - ▶ **Mercer's condition:** the kernels must satisfy Mercer's condition. Relevance vector regression To overcome the above limitations, Tipping introduces a novel model: Relevance Vector Regression

Source: Tipping, M.E., 2001. Sparse Bayesian learning and the relevance vector machine. *Journal of Machine Learning Research*, 1(Jun), pp.211-244

Limitations of SVR

- SVR has been used in many applications Nevertheless, they suffer some limitations:
 - ▶ **non-probabilistic:** SVR does not output probabilistic predictions;
 - ▶ **C and ϵ :** parameters of C and ϵ must be determined by cross-validation; and
 - ▶ **Mercer's condition:** the kernels must satisfy Mercer's condition. Relevance vector regression To overcome the above limitations, Tipping introduces a novel model: Relevance Vector Regression

Source: Tipping, M.E., 2001. Sparse Bayesian learning and the relevance vector machine. *Journal of Machine Learning Research*, 1(Jun), pp.211-244

Limitations of SVR

- SVR has been used in many applications Nevertheless, they suffer some limitations:
 - ▶ **non-probabilistic:** SVR does not output probabilistic predictions;
 - ▶ **C and ϵ :** parameters of C and ϵ must be determined by cross-validation; and
 - ▶ **Mercer's condition:** the kernels must satisfy Mercer's condition. Relevance vector regression To overcome the above limitations, Tipping introduces a novel model: Relevance Vector Regression

Source: Tipping, M.E., 2001. Sparse Bayesian learning and the relevance vector machine. *Journal of Machine Learning Research*, 1(Jun), pp.211-244

Limitations of SVR

- SVR has been used in many applications Nevertheless, they suffer some limitations:
 - ▶ **non-probabilistic:** SVR does not output probabilistic predictions;
 - ▶ **C and ϵ :** parameters of C and ϵ must be determined by cross-validation; and
 - ▶ **Mercer's condition:** the kernels must satisfy Mercer's condition. Relevance vector regression To overcome the above limitations, Tipping introduces a novel model: Relevance Vector Regression

Source: Tipping, M.E., 2001. Sparse Bayesian learning and the relevance vector machine. *Journal of Machine Learning Research*, 1(Jun), pp.211-244

Relevance Vector Regression (RVR)

- Assume the relationship between $\{\mathbf{x}_n\}$ and $\{y_n\}_{n=1}^N$ is:

$$y_n = \mathbf{w}^T \phi(\mathbf{x}) + \epsilon_n,$$

where ϵ_n are i.i.d with the Gaussian distribution having mean-zero and variance σ^2 .

- Due to the assumption of independence of the t_n , the likelihood of the complete data set can be written as

$$p(\mathbf{y}|\mathbf{w}, \sigma^2) = (2\pi\sigma^2)^{-N/2} \exp \left\{ -\frac{1}{2\sigma^2} \|\mathbf{y} - \Phi\mathbf{w}\|^2 \right\}, \quad (4)$$

where $\mathbf{y} = (y_1 \dots y_N)^T$, $\mathbf{w} = (w_0 \dots w_N)^T$ and Φ is the $N \times (N + 1)$ 'design' matrix with $\Phi = [\phi(\mathbf{x}_1), \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_N)]^T$, wherein $\phi(\mathbf{x}_n) = [1, K(\mathbf{x}_n, \mathbf{x}_1), K(\mathbf{x}_n, \mathbf{x}_2), \dots, K(\mathbf{x}_n, \mathbf{x}_N)]^T$ and $K(.,.)$ is a kernel function.

RVR (cont'd)

- Tipping assumes a zero-mean Gaussian prior distribution over \mathbf{w} :

$$p(\mathbf{w}|\alpha) = \prod_{i=0}^N \mathcal{N}(w_i|0, \alpha_i^{-1}), \quad (5)$$

with α a vector of $N + 1$ hyperparameters.

- The suitable priors for α and β are Gamma distributions:

$$p(\alpha) = \prod_{i=0}^N \text{Gamma}(\alpha_i|a, b),$$
$$p(\beta) = \text{Gamma}(\beta|c, d),$$

with $\beta \equiv \sigma^{-2}$ and where $\text{Gamma}(\alpha|a, b) = \Gamma(a)^{-1} b^a \alpha^{a-1} e^{-b\alpha}$, with $\Gamma(a) = \int_0^\infty t^{a-1} e^{-t} dt$, the 'gamma function'.

RVR (cont'd)

- Tipping assumes a zero-mean Gaussian prior distribution over \mathbf{w} :

$$p(\mathbf{w}|\alpha) = \prod_{i=0}^N \mathcal{N}(w_i|0, \alpha_i^{-1}), \quad (5)$$

with α a vector of $N + 1$ hyperparameters.

- The suitable priors for α and β are Gamma distributions:

$$p(\alpha) = \prod_{i=0}^N \text{Gamma}(\alpha_i|a, b),$$
$$p(\beta) = \text{Gamma}(\beta|c, d),$$

with $\beta \equiv \sigma^{-2}$ and where $\text{Gamma}(\alpha|a, b) = \Gamma(a)^{-1} b^a \alpha^{a-1} e^{-b\alpha}$, with $\Gamma(a) = \int_0^\infty t^{a-1} e^{-t} dt$, the 'gamma function'.

Pros and Cons of RVR

- **Pros**

- ▶ The relevance vector machine or RVM (Tipping, 2001) is a probabilistic regression model under the Bayesian framework;
- ▶ compared to the SVR, the RVR results in a sparser model;
- ▶ The kernels do not need to satisfy Mercer's condition;

- **Con** RVR only allows regression from multivariate inputs to a univariate output variable, but not to multiple variables. The operation of many engineering systems is influenced by multiple variables. For example, current and voltage are indispensable for electrical systems.

Pros and Cons of RVR

- **Pros**

- ▶ The relevance vector machine or RVM (Tipping, 2001) is a probabilistic regression model under the Bayesian framework;
- ▶ compared to the SVR, the RVR results in a sparser model;
- ▶ The kernels do not need to satisfy Mercer's condition;

- **Con** RVR only allows regression from multivariate inputs to a univariate output variable, but not to multiple variables. The operation of many engineering systems is influenced by multiple variables. For example, current and voltage are indispensable for electrical systems.

Pros and Cons of RVR

- **Pros**

- ▶ The relevance vector machine or RVM (Tipping, 2001) is a probabilistic regression model under the Bayesian framework;
- ▶ compared to the SVR, the RVR results in a sparser model;
- ▶ The kernels do not need to satisfy Mercer's condition;

- **Con** RVR only allows regression from multivariate inputs to a univariate output variable, but not to multiple variables. The operation of many engineering systems is influenced by multiple variables. For example, current and voltage are indispensable for electrical systems.

Pros and Cons of RVR

- **Pros**

- ▶ The relevance vector machine or RVM (Tipping, 2001) is a probabilistic regression model under the Bayesian framework;
- ▶ compared to the SVR, the RVR results in a sparser model;
- ▶ The kernels do not need to satisfy Mercer's condition;

- **Con** RVR only allows regression from multivariate inputs to a univariate output variable, but not to multiple variables. The operation of many engineering systems is influenced by multiple variables. For example, current and voltage are indispensable for electrical systems.

Literature review and our methods

RVR has been used to predict the remaining useful life.

- **Literature review**

- ▶ Some authors combine a set of univariate RVR models to output multiple variables, and
- ▶ Some authors regard that the weight matrix is inducted by separating into a vector distribution;

- **Our methods:**

- ▶ We propose a multivariate RVR model (MRVR), in which the weight matrix is a matrix Gaussian distribution;
- ▶ The hyperparameters of the MRVR model are estimated by Nesterov's Accelerated Gradient (NAG) method to obtain numerical solutions

Literature review and our methods

RVR has been used to predict the remaining useful life.

- **Literature review**

- ▶ Some authors combine a set of univariate RVR models to output multiple variables, and
- ▶ Some authors regard that the weight matrix is inducted by separating into a vector distribution;

- **Our methods:**

- ▶ We propose a multivariate RVR model (MRVR), in which the weight matrix is a matrix Gaussian distribution;
- ▶ The hyperparameters of the MRVR model are estimated by Nesterov's Accelerated Gradient (NAG) method to obtain numerical solutions

Literature review and our methods

RVR has been used to predict the remaining useful life.

- **Literature review**

- ▶ Some authors combine a set of univariate RVR models to output multiple variables, and
- ▶ Some authors regard that the weight matrix is inducted by separating into a vector distribution;

- **Our methods:**

- ▶ We propose a multivariate RVR model (MRVR), in which the weight matrix is a matrix Gaussian distribution;
- ▶ The hyperparameters of the MRVR model are estimated by Nesterov's Accelerated Gradient (NAG) method to obtain numerical solutions

Literature review and our methods

RVR has been used to predict the remaining useful life.

- **Literature review**

- ▶ Some authors combine a set of univariate RVR models to output multiple variables, and
- ▶ Some authors regard that the weight matrix is inducted by separating into a vector distribution;

- **Our methods:**

- ▶ We propose a multivariate RVR model (MRVR), in which the weight matrix is a matrix Gaussian distribution;
- ▶ The hyperparameters of the MRVR model are estimated by Nesterov's Accelerated Gradient (NAG) method to obtain numerical solutions

Table of Contents

- 1 A brief introduction to Remaining Useful Life and Relevance Vector Regression
- 2 Multivariate Relevance Vector Regression**
- 3 Case study

Multivariate RVR

- An MRVR is proposed as following

$$\mathbf{x}_{n+l} = \mathbf{W} \phi(\mathbf{x}_n) + \epsilon, \quad (6)$$

where

- ▶ $\mathbf{x}_{n+l} = [x_{1,n+l}, \dots, x_{M,n+l}]^T \in \mathbb{R}^M$ is the l -step forward prediction vector, and $1 < n+l \leq N$;
- ▶ $\phi(\mathbf{x}_n) = [\mathbf{1}, \mathcal{K}(\mathbf{x}_n, \mathbf{x}_1), \dots, \mathcal{K}(\mathbf{x}_n, \mathbf{x}_N)] \in \mathbb{R}^{N+1}$ denotes a design vector, in which $\mathcal{K}(\mathbf{x}_n, \mathbf{x}_j) \in \mathbb{R}$ is a kernel function,
- ▶ $\mathbf{W} \in \mathbb{R}^{M \times (N+1)}$ is a weight matrix of the design vector $\phi(\mathbf{x}_n) \triangleq \phi$, and
- ▶ ϵ is assumed to be a Gaussian distributed random error vector with the zero mean and a diagonal covariance matrix $\Sigma_0 = \text{diag}\{\sigma_1^2, \sigma_2^2, \dots, \sigma_M^2\} \in \mathbb{R}^{M \times M}$, and $\text{diag}(\cdot)$ denotes a diagonal matrix.

Multivariate RVR

- An MRVR is proposed as following

$$\mathbf{x}_{n+l} = \mathbf{W} \phi(\mathbf{x}_n) + \epsilon, \quad (6)$$

where

- ▶ $\mathbf{x}_{n+l} = [x_{1,n+l}, \dots, x_{M,n+l}]^T \in \mathbb{R}^M$ is the l -step forward prediction vector, and $1 < n+l \leq N$;
- ▶ $\phi(\mathbf{x}_n) = [\mathbf{1}, \mathcal{K}(\mathbf{x}_n, \mathbf{x}_1), \dots, \mathcal{K}(\mathbf{x}_n, \mathbf{x}_N)] \in \mathbb{R}^{N+1}$ denotes a design vector, in which $\mathcal{K}(\mathbf{x}_n, \mathbf{x}_j) \in \mathbb{R}$ is a kernel function,
- ▶ $\mathbf{W} \in \mathbb{R}^{M \times (N+1)}$ is a weight matrix of the design vector $\phi(\mathbf{x}_n) \triangleq \phi$, and
- ▶ ϵ is assumed to be a Gaussian distributed random error vector with the zero mean and a diagonal covariance matrix $\Sigma_0 = \text{diag}\{\sigma_1^2, \sigma_2^2, \dots, \sigma_M^2\} \in \mathbb{R}^{M \times M}$, and $\text{diag}(\cdot)$ denotes a diagonal matrix.

Multivariate RVR

- An MRVR is proposed as following

$$\mathbf{x}_{n+l} = \mathbf{W} \phi(\mathbf{x}_n) + \epsilon, \quad (6)$$

where

- ▶ $\mathbf{x}_{n+l} = [x_{1,n+l}, \dots, x_{M,n+l}]^T \in \mathbb{R}^M$ is the l -step forward prediction vector, and $1 < n+l \leq N$;
- ▶ $\phi(\mathbf{x}_n) = [\mathbf{1}, \mathcal{K}(\mathbf{x}_n, \mathbf{x}_1), \dots, \mathcal{K}(\mathbf{x}_n, \mathbf{x}_N)] \in \mathbb{R}^{N+1}$ denotes a design vector, in which $\mathcal{K}(\mathbf{x}_n, \mathbf{x}_j) \in \mathbb{R}$ is a kernel function,
- ▶ $\mathbf{W} \in \mathbb{R}^{M \times (N+1)}$ is a weight matrix of the design vector $\phi(\mathbf{x}_n) \triangleq \phi$, and
- ▶ ϵ is assumed to be a Gaussian distributed random error vector with the zero mean and a diagonal covariance matrix $\Sigma_0 = \text{diag}\{\sigma_1^2, \sigma_2^2, \dots, \sigma_M^2\} \in \mathbb{R}^{M \times M}$, and $\text{diag}(\cdot)$ denotes a diagonal matrix.

Multivariate RVR

- An MRVR is proposed as following

$$\mathbf{x}_{n+l} = \mathbf{W} \phi(\mathbf{x}_n) + \epsilon, \quad (6)$$

where

- ▶ $\mathbf{x}_{n+l} = [x_{1,n+l}, \dots, x_{M,n+l}]^T \in \mathbb{R}^M$ is the l -step forward prediction vector, and $1 < n+l \leq N$;
- ▶ $\phi(\mathbf{x}_n) = [\mathbf{1}, \mathcal{K}(\mathbf{x}_n, \mathbf{x}_1), \dots, \mathcal{K}(\mathbf{x}_n, \mathbf{x}_N)] \in \mathbb{R}^{N+1}$ denotes a design vector, in which $\mathcal{K}(\mathbf{x}_n, \mathbf{x}_j) \in \mathbb{R}$ is a kernel function,
- ▶ $\mathbf{W} \in \mathbb{R}^{M \times (N+1)}$ is a weight matrix of the design vector $\phi(\mathbf{x}_n) \triangleq \phi$, and
- ▶ ϵ is assumed to be a Gaussian distributed random error vector with the zero mean and a diagonal covariance matrix $\Sigma_0 = \text{diag}\{\sigma_1^2, \sigma_2^2, \dots, \sigma_M^2\} \in \mathbb{R}^{M \times M}$, and $\text{diag}(\cdot)$ denotes a diagonal matrix.

PDF of \mathbf{x}_{n+1}

Probability Density Function (PDF) of \mathbf{x}_{n+1} conditioned on \mathbf{W} and Σ_0 can be written by

$$p(\mathbf{x}_{n+1} | \mathbf{W}, \Sigma_0) = (2\pi)^{-\frac{M}{2}} |\Sigma_0|^{-\frac{1}{2}} \times \exp\left(-\frac{1}{2}(\mathbf{x}_{n+1} - \mathbf{W}\phi)^T \Sigma_0^{-1}(\mathbf{x}_{n+1} - \mathbf{W}\phi)\right), \quad (7)$$

here $|\cdot|$ is the determinant of a square matrix.

Prior and Posterior distributions of \mathbf{W}

- **Prior distributions of \mathbf{W}** To avoid the over-fitting problem of model (6), a prior matrix Gaussian distribution is assigned on the $M \times (N + 1)$ dimension weight matrix \mathbf{W} , which is denoted as $\mathbf{W} \sim \mathcal{MN}_{M,N+1}(\mathbf{0}, \mathbf{\Psi}, \mathbf{\Gamma})$, which gives

$$\begin{aligned} p(\mathbf{W} | \mathbf{\Psi}, \mathbf{\Gamma}) &= (2\pi)^{-\frac{M(N+1)}{2}} |\mathbf{\Psi}|^{-\frac{N+1}{2}} |\mathbf{\Gamma}|^{-\frac{M}{2}} \\ &\times \text{etr} \left(-\frac{1}{2} \mathbf{\Gamma}^{-1} \mathbf{W}^T \mathbf{\Psi}^{-1} \mathbf{W} \right), \end{aligned} \quad (8)$$

where $\text{etr}(\cdot)$ is the exponential of trace of a function of the trace of the matrix.

- **Posterior distributions of \mathbf{W}** The posterior distribution of weight matrix \mathbf{W} is matrix Gaussian, and its PDF is formulated in the following form.

$$\begin{aligned} p(\mathbf{W} | \mathbf{x}_{n+l}, \mathbf{\Psi}, \mathbf{\Gamma}, \mathbf{\Sigma}_0) &= (2\pi)^{-\frac{M(N+1)}{2}} |\tilde{\mathbf{\Psi}}|^{-\frac{N+1}{2}} \\ &\times |\tilde{\mathbf{\Gamma}}|^{-\frac{M}{2}} \text{etr} \left(-\frac{1}{2} \tilde{\mathbf{\Gamma}}^{-1} (\mathbf{W} - \tilde{\boldsymbol{\mu}})^T \tilde{\mathbf{\Psi}}^{-1} (\mathbf{W} - \tilde{\boldsymbol{\mu}}) \right). \end{aligned} \quad (9)$$

Prior and Posterior distributions of \mathbf{W}

- **Prior distributions of \mathbf{W}** To avoid the over-fitting problem of model (6), a prior matrix Gaussian distribution is assigned on the $M \times (N + 1)$ dimension weight matrix \mathbf{W} , which is denoted as $\mathbf{W} \sim \mathcal{MN}_{M,N+1}(\mathbf{0}, \mathbf{\Psi}, \mathbf{\Gamma})$, which gives

$$\begin{aligned} p(\mathbf{W} | \mathbf{\Psi}, \mathbf{\Gamma}) &= (2\pi)^{-\frac{M(N+1)}{2}} |\mathbf{\Psi}|^{-\frac{N+1}{2}} |\mathbf{\Gamma}|^{-\frac{M}{2}} \\ &\times \text{etr} \left(-\frac{1}{2} \mathbf{\Gamma}^{-1} \mathbf{W}^T \mathbf{\Psi}^{-1} \mathbf{W} \right), \end{aligned} \quad (8)$$

where $\text{etr}(\cdot)$ is the exponential of trace of a function of the trace of the matrix.

- **Posterior distributions of \mathbf{W}** The posterior distribution of weight matrix \mathbf{W} is matrix Gaussian, and its PDF is formulated in the following form.

$$\begin{aligned} p(\mathbf{W} | \mathbf{x}_{n+l}, \mathbf{\Psi}, \mathbf{\Gamma}, \mathbf{\Sigma}_0) &= (2\pi)^{-\frac{M(N+1)}{2}} |\tilde{\mathbf{\Psi}}|^{-\frac{N+1}{2}} \\ &\times |\tilde{\mathbf{\Gamma}}|^{-\frac{M}{2}} \text{etr} \left(-\frac{1}{2} \tilde{\mathbf{\Gamma}}^{-1} (\mathbf{W} - \tilde{\boldsymbol{\mu}})^T \tilde{\mathbf{\Psi}}^{-1} (\mathbf{W} - \tilde{\boldsymbol{\mu}}) \right). \end{aligned} \quad (9)$$

Parameter Estimation

- **Prediction distribution** The distribution of the predicted \mathbf{x}_{k+l} , based on the former prediction \mathbf{x}_{n+l} , can be obtained by

$$\begin{aligned} p(\mathbf{x}_{k+l}|\mathbf{x}_{n+l}) &= \iiint p(\mathbf{x}_{k+l}|\mathbf{W}, \Sigma_0) p(\mathbf{W}|\mathbf{x}_{n+l}, \Psi, \Gamma, \Sigma_0) \\ &\quad \times p(\Psi, \Gamma, \Sigma_0|\mathbf{x}_{n+l}) d\mathbf{W} d\Psi d\Gamma d\Sigma_0. \end{aligned} \quad (10)$$

- **Marginal likelihood function** The marginal likelihood function $p(\mathbf{x}_{n+l}|\Psi, \Gamma, \Sigma_0)$ can be obtained by integrating over the weight parameters \mathbf{W} as

$$\begin{aligned} p(\mathbf{x}_{n+l}|\Psi, \Gamma, \Sigma_0) &= \int p(\mathbf{x}_{n+l}|\mathbf{W}, \Sigma_0) p(\mathbf{W}|\Psi, \Gamma) d\mathbf{W} \\ &= \int p(\mathbf{x}_{n+l}|\text{vec}(\mathbf{W}^T), \Sigma_0) p(\text{vec}(\mathbf{W}^T)|\Psi, \Gamma) d\text{vec}(\mathbf{W}^T). \end{aligned} \quad (11)$$

Parameter Estimation

- **Prediction distribution** The distribution of the predicted \mathbf{x}_{k+l} , based on the former prediction \mathbf{x}_{n+l} , can be obtained by

$$\begin{aligned} p(\mathbf{x}_{k+l}|\mathbf{x}_{n+l}) &= \iiint p(\mathbf{x}_{k+l}|\mathbf{W}, \Sigma_0) p(\mathbf{W}|\mathbf{x}_{n+l}, \Psi, \Gamma, \Sigma_0) \\ &\quad \times p(\Psi, \Gamma, \Sigma_0|\mathbf{x}_{n+l}) d\mathbf{W} d\Psi d\Gamma d\Sigma_0. \end{aligned} \quad (10)$$

- **Marginal likelihood function** The marginal likelihood function $p(\mathbf{x}_{n+l}|\Psi, \Gamma, \Sigma_0)$ can be obtained by integrating over the weight parameters \mathbf{W} as

$$\begin{aligned} p(\mathbf{x}_{n+l}|\Psi, \Gamma, \Sigma_0) &= \int p(\mathbf{x}_{n+l}|\mathbf{W}, \Sigma_0) p(\mathbf{W}|\Psi, \Gamma) d\mathbf{W} \\ &= \int p(\mathbf{x}_{n+l}|\text{vec}(\mathbf{W}^T), \Sigma_0) p(\text{vec}(\mathbf{W}^T)|\Psi, \Gamma) d\text{vec}(\mathbf{W}^T). \end{aligned} \quad (11)$$

Parameter Estimation (cont'd)

- The negative log of the marginal likelihood is acquired

$$\begin{aligned} \mathcal{L}(\mathbf{x}_{n+l} | \Psi, \Gamma, \Sigma_0) &= \frac{1}{2} \ln |\Sigma_0| + \frac{N+1}{2} \ln |\Psi| + \frac{M}{2} \ln |\Gamma| \\ &+ \frac{1}{2} \ln |\mathbf{A}| + E(\boldsymbol{\mu}) + \frac{M}{2} \ln(2\pi), \end{aligned} \quad (12)$$

- set $\frac{\partial \mathcal{L}}{\partial \Psi_i} = 0$, $\frac{\partial \mathcal{L}}{\partial \Gamma_j} = 0$, and $\frac{\partial \mathcal{L}}{\partial \sigma_i^2} = 0$, it is difficult to obtain explicit solutions of the hyperparameters Ψ_i , Γ_j and σ_i^2 .
- The NAG (Nesterov's Accelerated Gradient) method is used to obtain numerical solutions of the hyperparameters

Parameter Estimation (cont'd)

- The negative log of the marginal likelihood is acquired

$$\begin{aligned} \mathcal{L}(\mathbf{x}_{n+l} | \Psi, \Gamma, \Sigma_0) &= \frac{1}{2} \ln |\Sigma_0| + \frac{N+1}{2} \ln |\Psi| + \frac{M}{2} \ln |\Gamma| \\ &+ \frac{1}{2} \ln |\mathbf{A}| + E(\boldsymbol{\mu}) + \frac{M}{2} \ln(2\pi), \end{aligned} \quad (12)$$

- set $\frac{\partial \mathcal{L}}{\partial \Psi_i} = 0$, $\frac{\partial \mathcal{L}}{\partial \Gamma_j} = 0$, and $\frac{\partial \mathcal{L}}{\partial \sigma_i^2} = 0$, it is difficult to obtain explicit solutions of the hyperparameters Ψ_i , Γ_j and σ_i^2 .
- The NAG (Nesterov's Accelerated Gradient) method is used to obtain numerical solutions of the hyperparameters

Parameter Estimation (cont'd)

- The negative log of the marginal likelihood is acquired

$$\begin{aligned} \mathcal{L}(\mathbf{x}_{n+l} | \Psi, \Gamma, \Sigma_0) &= \frac{1}{2} \ln |\Sigma_0| + \frac{N+1}{2} \ln |\Psi| + \frac{M}{2} \ln |\Gamma| \\ &+ \frac{1}{2} \ln |\mathbf{A}| + E(\boldsymbol{\mu}) + \frac{M}{2} \ln(2\pi), \end{aligned} \quad (12)$$

- set $\frac{\partial \mathcal{L}}{\partial \Psi_i} = 0$, $\frac{\partial \mathcal{L}}{\partial \Gamma_j} = 0$, and $\frac{\partial \mathcal{L}}{\partial \sigma_i^2} = 0$, it is difficult to obtain explicit solutions of the hyperparameters Ψ_i , Γ_j and σ_i^2 .
- The NAG (Nesterov's Accelerated Gradient) method is used to obtain numerical solutions of the hyperparameters

First hitting time

- The PDF of the degradation prediction takes the form

$$p(\mathbf{x}_{k+l}|\mathbf{x}_{n+l}) = \mathcal{N}(\mathbf{x}_{k+l}|\Sigma_0(\Sigma_0^{-1} \otimes \phi^T(\mathbf{x}_k))\Lambda\Sigma^{-1}\boldsymbol{\mu}, \Sigma_k), \quad (13)$$

- Given the observed measurements $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$, the RUL of each variable at time t_k is defined by

$$L_{ik} \triangleq L_i(t_k) = \inf\{t_l : x_i(t_k + t_l) \in \mathcal{B}_i | \mathbf{x}_{1:k}\}, \quad (14)$$

where $\inf\{\cdot\}$ denotes the infimum of a discrete; t_l represents the time length of the multi-step prediction; $x_i(t_k + t_l)$ is the degradation path at time $t_k + t_l$, $\mathbf{x}_{1:k}$ denotes the historical measurements from t_1 to t_k ; and \mathcal{B}_i refers to a boundary set, containing a boundary, barrier, or failure threshold.

First hitting time

- The PDF of the degradation prediction takes the form

$$p(\mathbf{x}_{k+l}|\mathbf{x}_{n+l}) = \mathcal{N}(\mathbf{x}_{k+l}|\Sigma_0(\Sigma_0^{-1} \otimes \phi^T(\mathbf{x}_k))\Lambda\Sigma^{-1}\boldsymbol{\mu}, \Sigma_k), \quad (13)$$

- Given the observed measurements $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k$, the RUL of each variable at time t_k is defined by

$$L_{ik} \triangleq L_i(t_k) = \inf\{t_l : x_i(t_k + t_l) \in \mathcal{B}_i | \mathbf{x}_{1:k}\}, \quad (14)$$

where $\inf\{\cdot\}$ denotes the infimum of a discrete; t_l represents the time length of the multi-step prediction; $x_i(t_k + t_l)$ is the degradation path at time $t_k + t_l$, $\mathbf{x}_{1:k}$ denotes the historical measurements from t_1 to t_k ; and \mathcal{B}_i refers to a boundary set, containing a boundary, barrier, or failure threshold.

RUL prediction

- the mean of the RUL is obtained by

$$E_i(t_k) = \sum_{L_{ik}=0}^{+\infty} L_{ik} \cdot p_i(L_{ik}), \quad (15)$$

where $p_i(L_{ik}) = \frac{\phi(\mathbf{g}_{i,k+1})\Delta\mathbf{g}_{i,k+1}}{1-\Phi(\mathbf{g}_{i,k})}$, $\phi(\cdot)$ and $\Phi(\cdot)$ are the PDF and cumulative distribution function of a standard normal random variable, respectively;

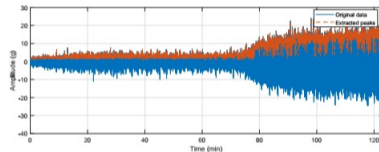
- the confidence interval for the prediction can also be obtained.

Table of Contents

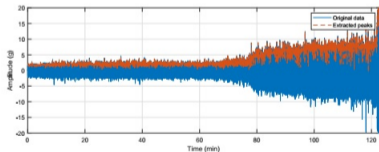
- 1 A brief introduction to Remaining Useful Life and Relevance Vector Regression
- 2 Multivariate Relevance Vector Regression
- 3 Case study**

Case study: Data

- A bearing dataset is used to demonstrate our proposed approach. Two accelerometers are placed on the bearings and positioned at 90° to each other, i.e., one is placed on the vertical axis and the other one is placed on the horizontal axis.
- Data from 78 minutes onwards are used for modelling.
- The operation is stopped when the amplitudes of the horizontal and vertical vibration signals are higher than $25g$ and $15g$, respectively.



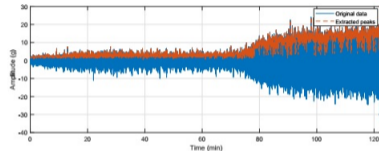
(a) The horizontal vibration signals.



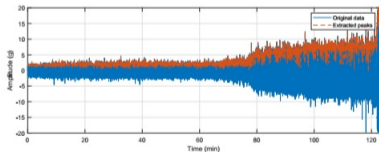
(b) The vertical vibration signals.

Case study: Data

- A bearing dataset is used to demonstrate our proposed approach. Two accelerometers are placed on the bearings and positioned at 90° to each other, i.e., one is placed on the vertical axis and the other one is placed on the horizontal axis.
- Data from 78 minutes onwards are used for modelling.
- The operation is stopped when the amplitudes of the horizontal and vertical vibration signals are higher than $25g$ and $15g$, respectively.



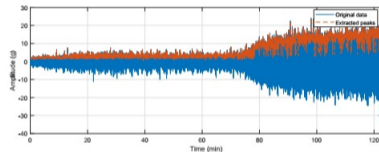
(a) The horizontal vibration signals.



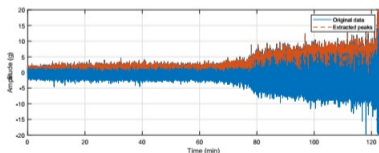
(b) The vertical vibration signals.

Case study: Data

- A bearing dataset is used to demonstrate our proposed approach. Two accelerometers are placed on the bearings and positioned at 90° to each other, i.e., one is placed on the vertical axis and the other one is placed on the horizontal axis.
- Data from 78 minutes onwards are used for modelling.
- The operation is stopped when the amplitudes of the horizontal and vertical vibration signals are higher than $25g$ and $15g$, respectively.



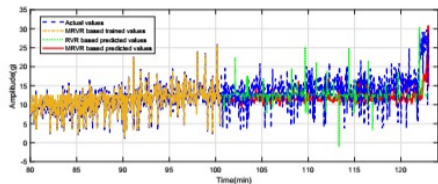
(a) The horizontal vibration signals.



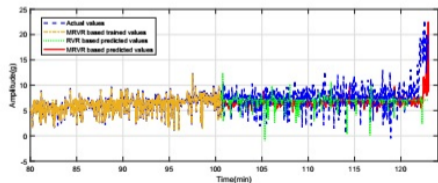
(b) The vertical vibration signals.

Case study: Degradation Path Prediction

- The peaks during 80 – 100min are used as inputs and those during 80.6 – 100.6min are used as outputs
- For the sake of comparison, both MRVR and RVR are built on the data
- the predicted degradation path based on the RVR cannot follow the actual vertical amplitude so well as that based on the MRVR.
- **Observation:** The MRVR outperforms the RVR.



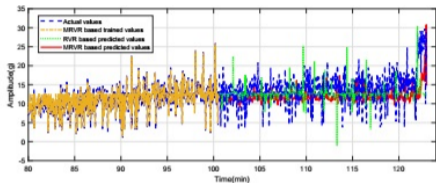
(a) The predicted amplitudes of the horizontal signal based on the RVR and MRVR.



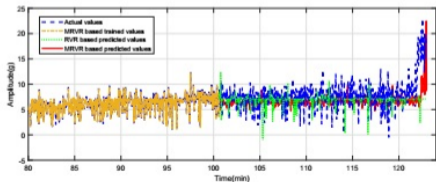
(b) The predicted amplitudes of the vertical signal based on the RVR and MRVR.

Case study: Degradation Path Prediction

- The peaks during 80 – 100min are used as inputs and those during 80.6 – 100.6min are used as outputs
- For the sake of comparison, both MRVR and RVR are built on the data
- the predicted degradation path based on the RVR cannot follow the actual vertical amplitude so well as that based on the MRVR.
- **Observation:** The MRVR outperforms the RVR.



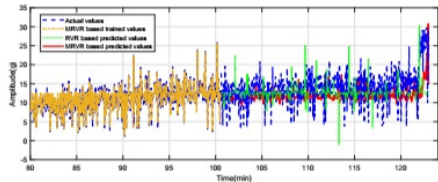
(a) The predicted amplitudes of the horizontal signal based on the RVR and MRVR.



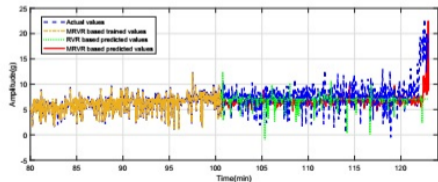
(b) The predicted amplitudes of the vertical signal based on the RVR and MRVR.

Case study: Degradation Path Prediction

- The peaks during 80 – 100min are used as inputs and those during 80.6 – 100.6min are used as outputs
- For the sake of comparison, both MRVR and RVR are built on the data
- the predicted degradation path based on the RVR cannot follow the actual vertical amplitude so well as that based on the MRVR.
- **Observation:** The MRVR outperforms the RVR.



(a) The predicted amplitudes of the horizontal signal based on the RVR and MRVR.



(b) The predicted amplitudes of the vertical signal based on the RVR and MRVR.

RUL prediction and Performance metrics

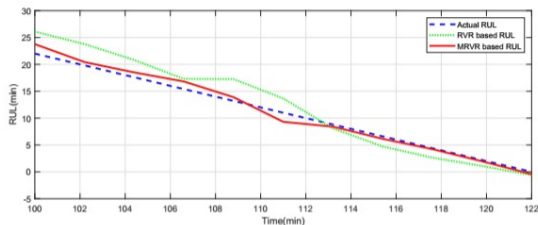


Figure: RUL Comparison

Performance metric	MRVR	RVR
MAE(min)	0.7960	2.3594
NRMSE(%)	8.7793	24.3193

Figure: Performance comparison

MAE=mean absolute error; NRMSE=Normalized Root Mean Relative Error

Conclusions

The MRVR outperforms!

Thank you all!

Questions?