

REMAINING USEFUL LIFE PREDICTION BASED ON A JOINT MODEL WITH DEGRADATION-FAILURE ASSOCIATION STRUCTURES

Xin HU¹, Xinbo QIAN², Xiao YANG³

Since the failures depend not only on internal degradation processes but also on external working conditions, failure thresholds of the performance indicator are stochastic for failure events. To improve remaining useful life (RUL) prediction accuracy, it is necessary to integrate both failure events and monitoring data, such as covariate-based hazard models. For most of the covariate-based hazard modeling methods, they essentially have a two-stage framework, degradation modeling first and then hazard modeling. However, the current two-stage method may ignore the influence of hazard on the degradation process, which may lead to significant bias in RUL prediction. A joint model is proposed to improve the RUL prediction performance by identifying the potential association structures between degradation and failures. The engine case study shows that the prediction performance of the is better than the two-stage method. Moreover, the effectiveness of the proposed method is reinforced by identifying the optimal association structure between degradation and failure.

Keywords: remaining useful life prediction, failure event, degradation, joint model, association structure identification

1. Introduction

Remaining useful life (RUL) prediction is one of the most important stages to prevent catastrophic failures in industrial systems. Accurate RUL prediction will effectively contribute to preventing unnecessary system

¹ Key Laboratory of Metallurgical Equipment and Control Technology, Ministry of Education, Wuhan University of Science and Technology, Wuhan, China, e-mail: shiwuxinya@163.com

² Corresponding author: Associate Professor, Hubei Key Laboratory of Mechanical Transmission and Manufacturing Engineering, Wuhan University of Science and Technology, Wuhan, China, e-mail: xinboqian@wust.edu.cn

³ Precision Manufacturing Institute, Wuhan University of Science and Technology, Wuhan, China

unavailability and massive downtime losses [1-2]. The commonly used methods to predict RUL are: physical-based, data-driven, and hybrid method. Hybrid method is one of the important methods for degradation prediction [3].

When only historical failure time data is available, the reliability model can be used for RUL prediction. However, when the failure time data is lacking, this reliability-based method does not perform well [4]. As the development of data acquisition techniques, more condition monitoring data can be available. Moreover, a feasible method is to apply the degradation models to RUL prediction [5]. Specifically, based on a large amount of condition monitoring (CM) data, the system degradation signal which is highly related to health status can be obtained. Such as the light intensity of the Light Emitting Diode and the resistance of the battery, which are commonly referred to in engineering as degradation signals of components. Evolution of these signals may lead to deterioration and final fail of component operation [6]. In the existing literature, a great deal of research work focused on the prediction of RUL by using observed data. Degradation is traditionally considered as a measured performance characteristic of cumulative changes over time leading to system failures. Moreover, many studies assume a constant failure threshold beyond which the degradation of the system will fail [7].

However, it is difficult, if not impossible, to predetermine such failure thresholds for devices with high-dimensional monitoring data. A typical example is the failure behavior of turbofan engine in aircraft system. Figure 1 shows lifetime and degradation data of turbofan engine [8], each broken line represents the degradation path of a specific engine, and the end of the straight line indicates engine failure. The data set was generated by commercial modular simulation software simulator (C-MAPSS) developed by NASA [9]. The lifetime of the degradation process of turbofan engines can be significantly influenced by the LPT coolant bleed factor. It is clear that due to external conditions, each engine fails at a different degradation level, and the influence of LPT coolant bleed rate on the lifetime is variable. Therefore, it may be difficult to define a specific failure threshold in advance. As Lee and Whitmore [10], Liu [11], and Song [12] et al. , RUL with fixed thresholds may underestimate or overestimate the real lifetime, resulting in additional costs or unexpected system failures. This fault behavior influenced by external conditions is not uncommon in practice.

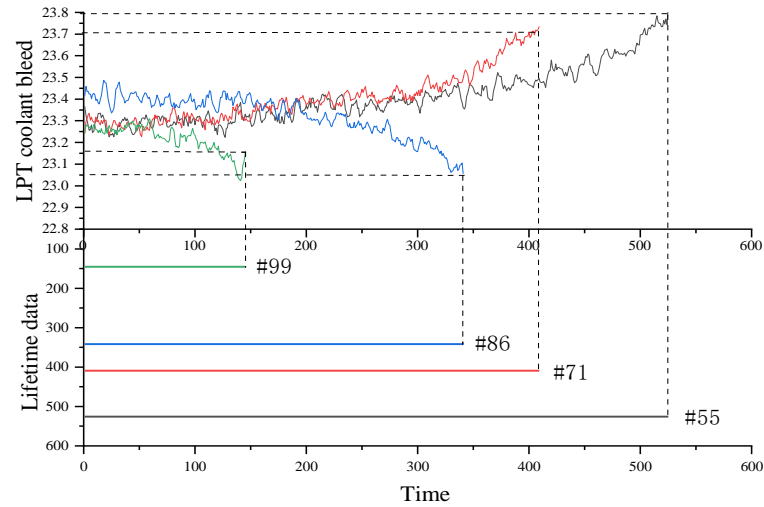


Fig. 1. Schematic diagram of random degradation threshold for failure events of aero-engines. The data comes from literature [8].

To predict RUL with stochastic thresholds, one possible approach is to assume that the failure thresholds follow a specified distribution [13]. However, the random threshold model lacks physical interpretation, so it may be difficult to correctly determine the fault threshold distribution for accurate RUL prediction. The literature [14] mentioned the poor prediction of RUL for lithium electronic batteries since the complexity of the electrochemical reactions inside the cell made modeling difficult. And it is difficult to collect data set based on the same operating conditions to characterize the degradation state of lithium batteries under real operating conditions [15]. Therefore, it is urgent to analyze failure time and degradation data together to improve the accuracy of the RUL prediction. Currently, it is popular to propose a two-stage method with degradation model and covariate-based hazard model for two stages respectively [16]. The degradation data can be treated as a time-varying covariate and substituted into the proportional hazards model for risk analysis. It is better to integrate the equipment service lifetime information with various state information. For example, Man et al [17] used the covariate-based hazard model for RUL prediction through simulation data, and adopted the two-stage method to estimate the parameters. According to the parameter estimation and degradation data, the conditional probability density function of in-service units was obtained, and then the RUL of in-service units was obtained. However, this method ignores the

monitoring error caused by repeated measurement and intermittent collection of monitoring values, and underestimates the association of model parameters. It needs to be improved when applied to industrial system RUL prediction [18].

Joint modeling of longitudinal and survival data is currently a popular framework in the medical field [19-20]. This modeling approach simultaneously analyzes repeated measurements and event outcomes, which can reduce bias in parameter estimation and improve the efficiency of statistical inference. However, there is relatively limited research and application in the field of reliability engineering. Moreover, there are many different characteristics in the degradation process, which will affect the accuracy of the covariate-based hazard model for RUL prediction. Such as the current value of the amount of degradation, the rate of degradation and the cumulative effect under the degradation trajectory [21]. To improve the accuracy of predicting RUL by the covariate-based hazard model, its core is to identify the association structure between degradation and failure. At present, the problem of selecting the most appropriate function form in a given data set has not been solved, and most of the work is focused on the process of the current degradation amount association failure time [22]. Therefore, this ignores the fact that different characteristics of the amount of degradation may also have an impact on the failure rate. The identification of potential association structures for degradation and failure is imminent.

In order to tackle this challenge, this study proposes a novel RUL method. Firstly, the failure events and degradation data are jointly modeled for RUL prediction. Then, the optimal potential association structure between degradation and failure is identified according to the prediction performance index of RUL. Specifically, the linear mixed model is applied to characterize system degradation. The potential degradation amount, degradation rate, and cumulative effect are included in the covariate-based hazard model as built-in covariates that affect the system failure rate. The model parameters are introduced into a Bayesian framework for simultaneous estimation, followed by RUL prediction. The optimal association structure for the joint model is identified based on the prediction performance of validation set. And the validity of the model is verified by test set. Moreover, this proposed method is applied to fit the data of turbofan engine in Figure 1 [8], [9], compared to the existing popular two-stage method. The main contributions of this work include:

- Compared with the existing two-stage model for modeling and analyzing degradation and failure data. The proposed joint model considers the influence of

failure rate on degradation process, and can estimate the parameters of degradation process and failure rate model simultaneously to correct the deviation.

- By identifying the fact that different characteristics of degradation have an impact on the failure rate, this significantly increases the model generalization for the industrial application.

The rest of the paper is organized as follows. Chapter II introduces the joint model with association structure of degradation and failure. Chapter III presents the association structure identification between degradation and failure. Chapter IV applies the proposed method to the Case Study. Chapter V concludes the study and discusses possible future works.

2. Joint model with association structure of degradation and failure

2.1 Joint modeling framework

This subsection focuses on the joint modeling framework with association structure identification between degradation and failure. The flowchart of the proposed method for RUL prediction with optimal association structure being identified, is as follows.

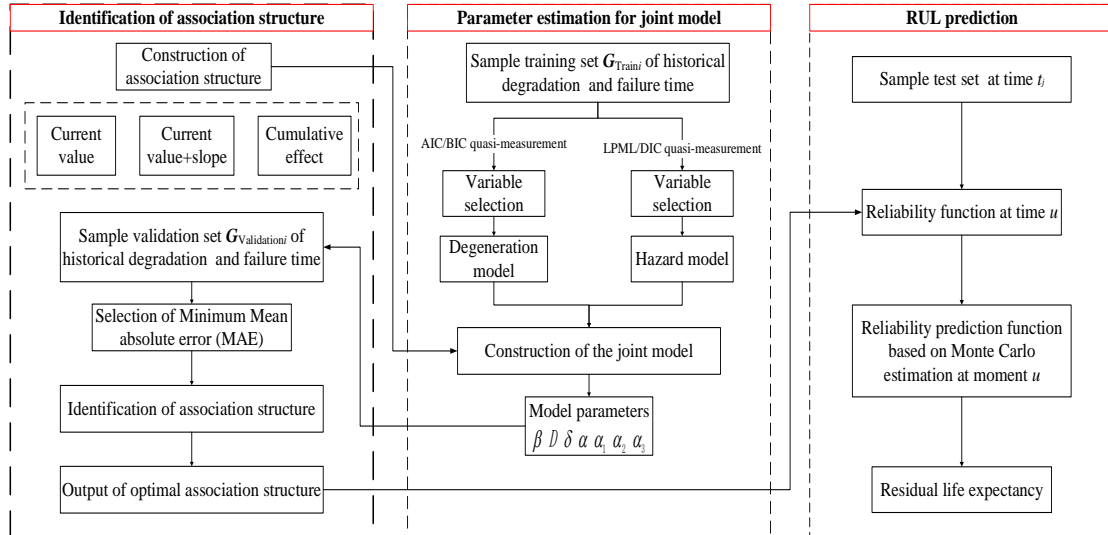


Fig. 2. Flowchart of the proposed method for RUL prediction by joint model with optimal association structure identification

The proposed method includes three main parts: parameter estimation of joint model, identification of association structure, and RUL prediction. The parameter estimation of the joint model introduces the random effect describing

the difference of samples, which serves as a joint basis for the degeneration model and the hazard model, and joint modeling of potential failure events and degradation data. It mainly includes model variable selection of the AIC/BIC and LPML/DIC criterion. And the posterior mean and posterior variance of the target parameters are estimated by random sampling with Markov Monte Carlo algorithm (MCMC) for Bayesian inference. The identification association structure is based on the joint modeling of degradation and failure functions under different association structures. Moreover, the optimal association structure is output based on the prediction performance evaluation index (MAE). The RUL predictions are mainly constructed by constructing the reliability function together with the test set samples and the relevant parameters estimated by the joint model. Then the reliability prediction function is obtained by updating the parameters through Monte Carlo estimation. Finally, the prediction result of RUL is obtained by integrating the reliability prediction function.

2.2 Degradation modeling and hazard modeling

Let $\{y_\lambda(t_{\lambda j}), t_{\lambda j}, \delta_\lambda; j = 1, 2, 3, \dots, n_\lambda\}$ denotes the data structure observed from n individuals. $y_\lambda(t_{\lambda j})$ denotes the monitoring data of the measured subject individual λ at the time points $t_{\lambda j}$ and $t_\lambda = \min(T_\lambda^*, C_\lambda)$ is the time of observation of the event of interest. It is assumed that the true event time T_λ^* and the truncation time C_λ are independent of each other. The joint model of the data consists of two sub-models defined by the failure events and the monitoring data [23]. The monitoring eigenvalues are modeled using a linear mixed model as follows

$$y_\lambda(t) = \eta_\lambda(t) + \varepsilon_\lambda(t) = X_\lambda^T(t)\beta + Z_\lambda^T(t)b_\lambda + \varepsilon_\lambda(t), \quad (1)$$

where $y_\lambda(t)$ denotes the time series of monitoring data at any time point t of the λ th individual, X_λ^T denotes the design matrix of the fixed effects β , and Z_λ^T is the design matrix of the random effect b_λ , where $b_\lambda \sim N(0, D)$, $\varepsilon_\lambda(t)$ is the measurement error, $\varepsilon_\lambda(t) \sim N(0, \sigma^2)$. The error terms $\varepsilon_\lambda(t)$ and the random effect b_λ are independent of each other, $\eta_\lambda(t)$ denotes the true value of the monitoring variable at time point t .

For the failure processes, a proportional hazards model is used to describe the risk of an event. Let the event time $T_\lambda = \min(T_\lambda^*, C_\lambda)$, T_λ^* represent the actual observed event time of the λ th individual and C_λ represent the truncation time. The form of the proportional hazards model is as follows

$$\begin{aligned}
h_\lambda(t | H_\lambda(t), m_\lambda) &= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} \Pr\{t \leq T_\lambda^* < t + \Delta t | T_\lambda^* \geq t, H_\lambda(t), m_\lambda\} \\
&= h_0(t) \exp\left[\gamma^* m_\lambda + f\{\eta_\lambda(t), \mathbf{b}_\lambda, \boldsymbol{\alpha}\}\right], \quad t > 0,
\end{aligned} \tag{2}$$

where $H_\lambda(t) = \eta_\lambda(s), 0 \leq s < t$ denotes the historical time series of potential monitoring values up to t , where $h_0(t)$ is the baseline risk function, m_λ is the vector of baseline covariates and the corresponding vector of regression coefficients is γ . The parameter vector α describes the strength of the association between monitoring eigenvalues and the event process, quantifying the impact of potential monitoring data $\eta_\lambda(t)$ on event risk during potential degradation. Let $\exp(\gamma_j)$ denotes the risk ratio for a unit change in $m_{\lambda j}$ at any time t , $\exp(\alpha_1)$ denotes the relative increase in survival risk at the same time $\exp(\alpha_1)$ for each unit increase in $\eta_\lambda(t)$ at time t . The various association structures of the functional form $f(\cdot)$ are described in detail in chapter 3. To complete the description of the failure risk process, we need to make appropriate assumptions about the baseline hazard function $h_0(t)$. To model this function while still considering flexibility, we use a penalized B-spline approximation for the baseline hazard. In particular, the logarithm of the baseline hazard function is expressed as

$$\log h_0(t) = \gamma_{h_0,0} + \sum_{q=1}^Q \gamma_{h_0,q} B_q(t, \mathbf{v}). \tag{3}$$

Here $B_q(t, \mathbf{v})$ denotes the q th basis function of the B spline with node v_1, \dots, v_Q and a vector of γ_{h_0} spline coefficients, increasing the number of nodes Q increases the flexibility of approximating $\log h_0(\cdot)$. However, we should balance the bias and variance to avoid overfitting. In the Bayesian framework, different association structures can be specified by targeting the form of function $f(\cdot)$. The detailed procedure is described in chapter 3.

In this paper, the deficit pool information criterion (AIC) and the Bayesian information criterion (BIC) are used to select the variables of the degradation model. Wang (2007) [24] pointed out that the selection of adjustment parameters may lead to overfitting and proposed the use of BIC for variable selection. The model with the smallest AIC and BIC is usually chosen when selecting parameters from a set of selected model variables. The bias information criterion (DIC) and the log pseudo-marginal likelihood (LPML) are used for hazard model variable selection. DIC is smaller indicate the better model fit, and the larger value of LPML indicates the better model fit.

2.3 Parameter estimation for joint model

To fully consider the potential relationships between the data, the potential

measurements of the degradation process, i.e., the true values without errors, and the failure time data are jointly modeled. The model parameters are estimated by using JMBayes package of the R software [25]. The identification is mainly based on Markov chain Monte Carlo (MCMC). Under the premise of given random effect, it is assumed that the degradation process and the failure process are independent, and the time series responses of each subject is independent. The likelihood function expression of the model parameters is derived as follows

$$p(y_\lambda, T_\lambda, \delta_\lambda | \mathbf{b}_\lambda, \phi) = p(y_\lambda | \mathbf{b}_\lambda, \phi) p(T_\lambda, \delta_\lambda | \mathbf{b}_\lambda, \phi), \quad (4)$$

$$p(y_{\lambda l} | \mathbf{b}_\lambda, \phi) = \prod_{\lambda} p(y_{\lambda l} | \mathbf{b}_\lambda, \phi), \quad (5)$$

where ϕ is the full parameter vector and $p(\cdot)$ is the appropriate probability density function. Under these assumptions, the posterior distribution is similar to

$$p(\phi, \mathbf{b}) \propto \prod_{\lambda=1}^n \prod_{l=1}^{n_\lambda} p(y_{\lambda l} | \mathbf{b}_\lambda, \phi) p(T_\lambda, \delta_\lambda | \mathbf{b}_\lambda, \phi) p(\mathbf{b}_\lambda, \phi) p(\phi), \quad (6)$$

$$p(y_{\lambda l} | \mathbf{b}_\lambda, \phi) = \exp \left\{ \left[y_{\lambda l} \psi_{\lambda l}(\mathbf{b}_\lambda) - c \{ \psi_{\lambda l}(\mathbf{b}_\lambda) \} \right] / a(\phi) - d(y_{\lambda l}, \phi) \right\}, \quad (7)$$

where $\phi^T = (\phi_t^T, \phi_y^T, \phi_b^T)$ denotes the complete parameter vector, ϕ_t denotes the parameters of the event time outcome, ϕ_y denotes the parameters of the degradation outcome, and ϕ_b denotes the unique parameters of the random effects covariance matrix. Formula (7) $\psi_{\lambda l}(\mathbf{b}_\lambda)$ and ϕ respectively represent the natural and dispersion parameters in the index family, and $c(\cdot)$, $a(\cdot)$, and $d(\cdot)$ are known functions that specify the members of the index family. For the survival function part

$$p(T_\lambda, \delta_\lambda | \mathbf{b}_\lambda, \phi) = h_\lambda(T_\lambda | H_\lambda(T_\lambda))^{\delta_\lambda} \exp \left\{ - \int_0^{T_\lambda} h_\lambda(s | H_\lambda(s)) ds \right\}, \quad (8)$$

$h_\lambda(\cdot)$ is given by formula (2)

$$S_\lambda(t | H_\lambda(t), \mathbf{m}_\lambda) = \exp \left\{ - \int_0^t h_0(s) \exp \left[\gamma^T \mathbf{m}_\lambda + f \{ \eta_\lambda(s), \alpha \} \right] ds \right\}. \quad (9)$$

For the parameter ϕ , this paper adopts the standard prior distribution. In particular, for the fixed effects vector of the degradation model β , the regression parameters of the survival model γ , the vector of spline coefficients of the baseline hazard γ_{h0} , and the association parameter α , which use an independent univariate diffusion normal prior. The joint likelihood function of the failure events and the degradation data integrates all the information of the two parts of the data and has a more complex structure. Therefore, it is difficult to obtain an analytic solution for the posterior $\pi(\phi | D)$. In this paper, MCMC simulation iterations based on the Metropolis Hasting sampling method is used to obtain Monte Carlo samples. Then, according to these Monte Carlo samples, the

posterior mean and posterior variance of parameters are estimated, and Bayesian inference is made.

3. Association structure identification between degradation and failure

In the joint framework, different association structures can be identified for the form of the function $f(\cdot)$ of formula (2) in chapter 2.2. In this paper, according to the reference [26], where also consider mainly three kinds of association structures in the failure process.

3.1 Association structure between current degradation value and failure

The current value indicates that the event failure rate at moment t is related to the degradation trajectory. The degradation and failure potential association structure are to establish the association between the potential measurement process in the degradation process, i.e., the measurement process that does not contain errors, with failure events. At this point, the joint model association term is $f(\cdot) = \alpha\eta_\lambda(t)$. The specific expression of the failure event is shown in formula (10)

$$h_\lambda(t) = h_0(t) \exp\{\gamma m_\lambda + \alpha\eta_\lambda(t)\}. \quad (10)$$

Here $h_0(\cdot)$ is the baseline risk function, m_λ is the vector of baseline covariates, and the corresponding vector of regression coefficients is γ . The parameter vector α describes the strength of the association between potential monitoring eigenvalues and failure time. And it quantifies the effect of the potential measurement process $\eta_\lambda(t)$ on the risk of failure during degradation.

3.2 Association structure between current value and slope of degradation and failure

The current value and slope indicates that the event failure rate at moment t is related to the degenerate trajectory and the slope of the degenerate trajectory at moment t . The potential association structure between degradation and failure events is to associate potential measurement processes in the degradation process, i.e., measurement processes that do not contain errors with failure events. At this point, the joint model association term is $f(\cdot) = \alpha_1\eta_\lambda(t) + \alpha_2\dot{\eta}_\lambda(t)$. The specific expression for the failure event is shown in formula (11)

$$h_\lambda(t) = h_0(t) \exp\{\gamma m_\lambda + \alpha_1\eta_\lambda(t) + \alpha_2\dot{\eta}_\lambda(t)\}, \quad \dot{\eta}_\lambda(t) = \frac{d\eta_\lambda(t)}{dt}. \quad (11)$$

Here $h_0(\cdot)$ is the baseline risk function, m_λ is the vector of baseline covariates, and the corresponding vector of regression coefficients is γ . The parameter vector α_1 describes the strength of the association between the potential monitoring eigenvalues and the failure time. The parameter vector α_2 describes the strength of the association between the slope of the monitoring trajectory and the failure time.

And it quantifies the influence of the potential measurement process $\eta_\lambda(t)$ and the slope $\dot{\eta}_\lambda(t)$ on the failure risk during the degradation process.

3.3 Association structure between cumulative effect of degradation and failure

The cumulative effect indicates that the event failure rate at moment t is related to the entire area under the degraded trajectory up to moment t . The potential association structure between degradation and failure is to associate potential measurement processes in the degradation process, i.e., measurement processes that do not contain errors with failure events. At this point, the joint model association term is $f(\cdot) = \alpha_3 \int_0^t \eta_\lambda(s) ds$. The specific expression for the failure event is shown in formula (12)

$$h_\lambda(t) = h_0(t) \exp\{\gamma m_\lambda + \alpha_3 \int_0^t \eta_\lambda(s) ds\}. \quad (12)$$

Here $h_0(\cdot)$ is the baseline risk function, m_λ is the vector of baseline covariates, and the corresponding vector of regression coefficients is γ . The parameter vector α_3 describes the strength of the association between the entire area under the potential degradation trace and the failure event. And it quantifies the influence of the entire area $\int_0^t \eta_\lambda(s) ds$ under the potential degradation trace on the failure risk during degradation.

3.4 Identification of association structure between degradation and failure by RUL prediction

To identify the optimal association structure, this chapter focuses on the RUL prediction based on the joint model of different association structures between degradation and failure, and then identifies it based on the predicted performance index (MAE). Please see the Appendix for details on how to do this. The joint model is fitted on a sample of size n , based on the monitoring data $y_\lambda(t) = \{y_\lambda(s); 0 \leq s \leq t\}$ of the new object i , indicating the degradation process up to the moment t degradation data. The RUL prediction is predicated on obtaining the conditional reliability of the sample test set. We pay more attention to the survival time $u > t$, and get the expression of the built reliability function as follows [27]

$$R_\lambda(u | t) = \Pr(T_\lambda^* \geq u | T_\lambda^* \geq t, y_\lambda(t), D_n), \quad (13)$$

the formula $D_n = \{T_\lambda, \delta_\lambda, y_\lambda; \lambda = 1, \dots, n\}$ represents the fitted joint model sample data, when the test sample $t' > t$ recorded new information can be updated prediction. According to this, the prediction is made according to formula (8), and formula (13) can be expressed as

$$R_\lambda(u | t) = \int \Pr(T_\lambda^* \geq u | T_\lambda^* \geq t, y_\lambda(t), \phi) p(\phi | D_n) d\phi, \quad (14)$$

the first part of the quilt product function is calculated by making full use of the conditional independence assumptions of formulas (4) and (5), and the first part of the quilt product function can be rewritten as

$$\begin{aligned} \Pr(T_\lambda^* \geq u | T_\lambda^* \geq t, y_\lambda(t), \phi) &= \int \Pr(T_\lambda^* \geq u | T_\lambda^* \geq t, b_\lambda(t), \phi) p(b_\lambda | T_\lambda^* \geq t, y_\lambda(t), \phi) db_\lambda \\ &= \int \frac{S_\lambda\{u | H_\lambda(u, b_\lambda), \phi\}}{S_\lambda\{t | H_\lambda(t, b_\lambda), \phi\}} p(b_\lambda | T_\lambda^* \geq t, y_\lambda(t), \phi) db_\lambda, \end{aligned} \quad (15)$$

where $S_\lambda(\cdot)$ is given by formula (9) and the degenerate historical data $H_\lambda(\cdot)$ is a function of random effects and parameters and is approximated by a linear mixed model. The first part of the quantile function, as shown above, is given by formula (15). In the second part, i.e., the posterior distribution of the parameters given the observed data, we use the standard asymptotic Bayesian approach with ϕ denoting the great likelihood estimate and H denoting the asymptotic variance matrix. And the posterior converges to a multivariate normal distribution $\{\phi | D_n\} \sim N(\phi, H)$ when the sample n is sufficiently large. The Metropolis-Hastings algorithm is combined with formulas (14) and (15) and multivariate t for Monte Carlo estimation, and the following results are obtained

$$R_\lambda(u | t) = L^{-1} \sum_{l=1}^L R_\lambda^{(l)}(u | t), \quad (16)$$

where L denotes the number of Monte Carlo samples, when $u > t$, the prediction of the RUL of the joint model at the i th sample at a time u in the future is denoted as

$$RUL_\lambda = \int_t^\infty R_\lambda(u | t), \quad (17)$$

the evaluation index of the predictive performance of the joint model, mean absolute error (MAE), is used to identify the degradation and failure potential association structure, and the optimal joint model of association structure between degradation and failure is output. The specific expression is as follows

$$MAE = K^{-1} \sum_{\lambda=1}^K |RUL_\lambda(t_j) - RUL_\lambda(t_j)|. \quad (18)$$

Here K denotes the sample capacity of the test set, and $RUL_\lambda(t_j)$ and $RUL_\lambda(t_j)$ denote the RUL predicted and true values of the λ th sample at the starting moment t_j of the prediction, respectively.

4. Case Study

4.1 Introduction to the case dataset

In this case study, the proposed joint model will be implemented and

evaluate based on the degradation data set of turbofan engines provided in reference [9]. The turbofan engine data is derived from the simulation software C-MAPSS and contains data on the life of the engine as it degrades to non-functional with increasing operating time. The monitoring data for each flight cycle consists of 26 dimensions of characteristic data (such as W32(LPT coolant flow)), 3 dimensions are the common flight conditions of the aircraft (containing flight altitude (OS_1), Mach numbers(OS_2) and temperature (OS_3)). For the original data information, First of all, the data set is divided, then noise reduction is applied. To eliminate the dimensional influence of data indicators, data standardization is required, and the Min-Max Normalization method is used.

To demonstrate the effectiveness of the proposed model, the authors adopt the Monte Carlo Cross Validation method for validation. The original data were randomly divided into training set G_{Train_i} , validation set $G_{Validation_i}$ and test set according to 88:11:1, i denotes the number of divisions, the sample segmentation was repeated $i=10$ times. In each segmentation, RUL prediction was performed for $k=11$ sample devices in the validation set under different prediction starting points, and k denotes the sample size of $G_{Validation_i}$. The prediction starting points were divided according to the percentage of the full lifetime of the samples. That is, $t_j=10\%,20\%,30\%.....90\%$. The average absolute error (MAE) of the prediction RUL for different $G_{Validation_i}$ is calculated, and the box line plots are drawn based on different prediction starting points and the prediction performance is compared with the two-stage method. Then, the association structure is identified based on the prediction performance criterion MAE, and the optimal joint model of association structure between degradation and failure is output. This specific method is provided in Appendix. The generalization ability of the optimal model is evaluated by test set.

4.2 Model variables selection

The hazard function significance analysis was performed on the monitored variables by R language, and the p-value results are shown in Table 1. Degradation model using the linear mixed effects model, the most significant variable W32 (LPT coolant flow) was selected as the dependent variable of the degradation model. Three flight conditions (flight altitude, Mach number, and temperature) that do not vary with time were considered as fixed effects. The redundant variables temperature is eliminated, and the interaction between time and altitude is considered. The differences in monitoring moments of different samples were considered as random effects, and monitoring errors are considered. Since the linear mixed effects model requires the assumption that the data set is to obey a normal distribution, a normality test is done for the W32 data set, and the

data meet the characteristic requirements of the linear mixed effects model requiring a normal distribution.

Table 1

P-value of the variable			
Parameter Type	OS_1	OS_2	W32
p-value	0.1662	0.0783	0.0230

A linear mixed effects model with all variables is

$$y_{\lambda} = \beta_0 + \beta_1 \cdot year + \beta_2 \cdot os_1 + \beta_3 \cdot os_2 + \beta_4 \cdot os_1 \cdot year + z_0 + z_1 \cdot year + \varepsilon_{\lambda}, \quad (19)$$

Data analysis using a linear mixed effects model yielded the following table of p-values for the coefficient estimates of each variable.

Table 2

P-values of the coefficient estimates for each variable					
Variables	(Intercept)	$year$	OS_1	OS_2	$OS_1 \cdot year$
p-value	0.0923	0.8932	0.5158	0.7663	0.9263

According to Table 2, we can analyze the p values of the four variables in formula (19) $year$, OS_1 , OS_2 , and $OS_1 \cdot year$ are insignificant. So the stepwise method is used to select the appropriate variables, and the results are shown in Table 3.

Table 3

Comparison of degradation models with different variable selection								
Variables	Formula (24)	Remove $year$	Remove OS_1	Remove OS_2	Remove $year OS_1$	Remove $year OS_2$	Remove $OS_1 OS_2$	Remove $year OS_1 OS_2$
AIC	-48870.2	-48880.2	-48870.6	-48874.6	-48880.9	-48884.6	-48875.1	-48885.4
BIC	-48799.3	-48817.2	24448.1	-48811.6	-48825.7	-48829.5	-48820.0	-48838.2
Loglik	24444.1	-48807.1	24443.3	24445.3	24447.5	24449.3	24444.6	24448.7

Preference is given to the linear mixed effects model with the minimum AIC and BIC value as the object of the final monitoring data for the selected model

$$y_{\lambda} = \beta_0 + \beta_4 \cdot os_1 \cdot year + z_0 + z_1 \cdot year + \varepsilon_{\lambda}. \quad (20)$$

According to the previous Table 1, the most significant variable affecting the hazard function is W32 (LPL cooling flow), followed by OS_2 . Similarly, OS_2 is excluded as a baseline covariate of the proportional hazards model, this method is

based on the joint LPML and DIC quasi-measurement analysis. The expressions are as follows

$$h_{\lambda}(t) = h_0(t) \exp \left[f \{ \eta_{\lambda}(t), \mathbf{b}_{\lambda}, \boldsymbol{\alpha} \} \right], \quad t > 0. \quad (21)$$

4.3 Joint model results of difficult association structures

The results of the joint model with association structure between current degradation value and failure (Table 4), Which reveal the potential degradation trajectory is significantly influenced by the strength of the association with the failure event process. The results of the joint model with association structure between current value and slope of degradation and failure are found (Table 4). The potential degradation trajectory has a significant effect on the strength of association with the failure event process, and the potential degradation rate has a non-significant effect on the strength of association with failure events. The results of the joint model with association structure between degradation and failure of cumulative effect (Table 4), which reveal that the cumulative effect of potential degradation trajectories is significantly influenced by the strength of the association with the failure event process. Thus, the RUL prediction performance of the joint Bayesian model based on different degradation and failure association structures is continued to be investigated. The results of the joint model construction are shown in the following table.

Table 4

Joint model results of different association structures between degradation and failure

	Variables	Coefficient	Standard Error	std.Dev	P value
Current value	Intercept	0.4665	7e-04	2e-03	<0.001
	$OS_{i,year}$	-0.0005	1e-04	1e-04	<0.001
	Assoct	-4.0318	0.0706	0.7724	<0.001
Current value+ slope	Intercept	0.4695	7e-04	0.0019	<0.001
	$OS_{i,year}$	-0.0008	1e-04	0.0001	<0.001
	Assoct	-3.8814	0.0843	0.8267	<0.001
	AssoctE	-0.0020	0.1551	3.1339	0.989
Cumulative effect	Intercept	0.471	4e-04	0.0016	<0.001
	$OS_{i,year}$	-0.0001	1e-04	0.0001	0.0905
	Assoct	-0.015	0.0009	0.0065	<0.001

4.4 Joint model with association structure optimization

In this paper, Monte Carlo Cross Validation method is used to randomly

divide train set G_{Train_i} and validation set $G_{Validation_i}$ (excluding test set) of the sample data. The percentage of the whole lifetime based on the test set data is used as the prediction starting point, calculate the performance index MAE predicted by test set RUL. The results of using different joint models of association structure between degradation and failure after cross-validation is presented in Figure 3.

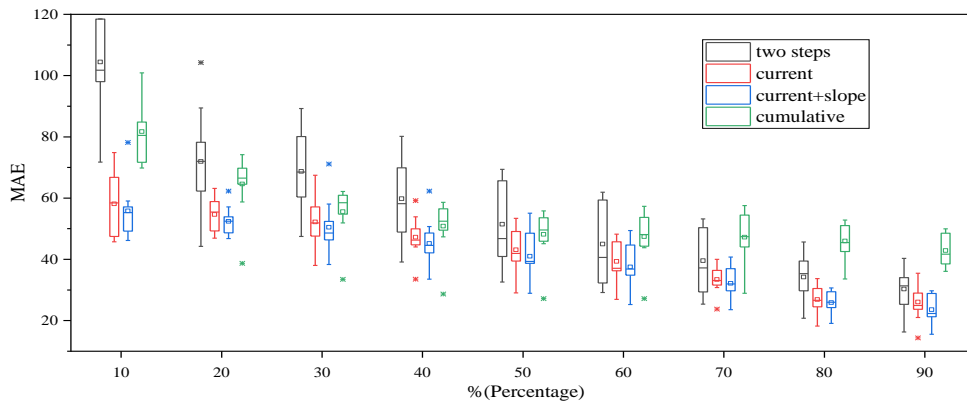


Fig. 3. The MAE of the RUL prediction from the Monte Carlo Cross Validation method

From the Figure 3, the prediction performance of the joint model of association structure between current degradation value and failure is significantly better than that of the traditional two-stage method. And it is especially significant at the early stage (i.e., the starting point of prediction is located at 10%-30% of the whole lifetime). The validity of the model is verified. Furthermore, regarding the identification of association structure, the joint model with association structure between current value and slope of degradation and failure has the best overall performance. And in the early stage, the middle stage (40%-60%) and the late stage (70%-90%). The second best prediction performance of RUL is the association structure between current degradation value and failure. The third predictive performance of RUL is the association structure between current degradation value and failure, which has a slow decreasing trend in the late stage. The prediction accuracy is improved with the increase of data, which verifies the validity of the model again.

To better illustrate the effectiveness of the proposed model in this paper, Figure 4 shows in detail the RUL prediction results of the joint model with different association structures between degradation and failure for individual #5 by test set.

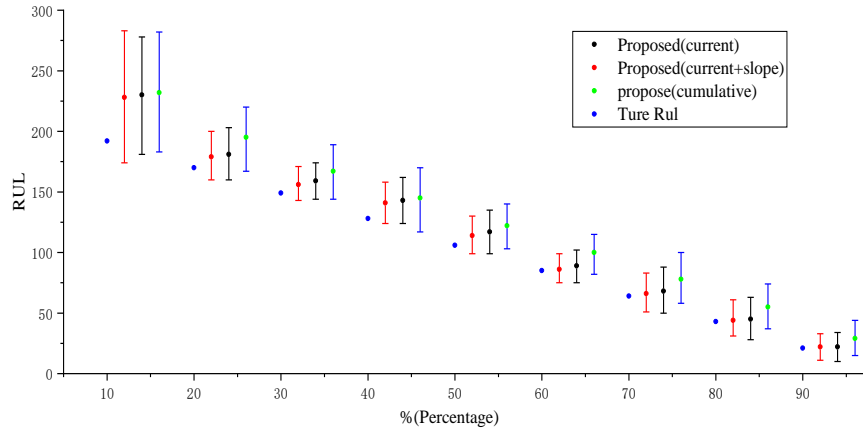
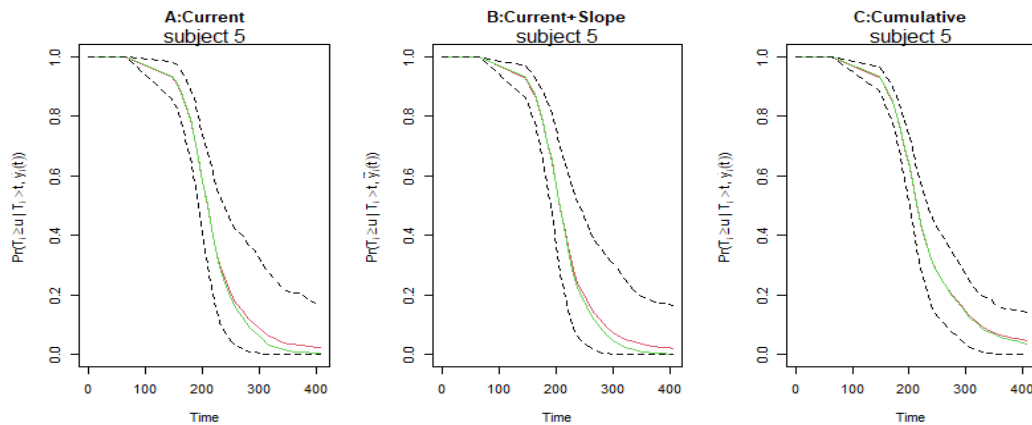


Fig. 4. The true and the estimated RUL of individual #5 at different prediction starting points via different association structures, the band shows the 95% prediction interval.

For individual #5, it is found that 95% of the prediction intervals contained the true RUL for all predicted moments t_j . Nevertheless, the RUL prediction is more important when the individual approaches its lifetime, the RUL prediction becomes much more accurate and approaches the true value when more data are collected. In the experimental results, it is verified that the joint model with association structure between current value and slope of degradation and failure has the best prediction performance. And in the early stage (i.e., the prediction starting point moment is located at 10%-30% of the whole lifetime), the middle stage (40%-60%) and the late stage (70%-90%). More detailed prediction results are provided for Individual #5 under the joint model with different association structures. The predicted starting time based on the estimated reliability function and point-wise quartiles is located at 30% and 80% of lifetime as shown in Figure 5.



(a) Prediction starting point $t_j = 30\%$

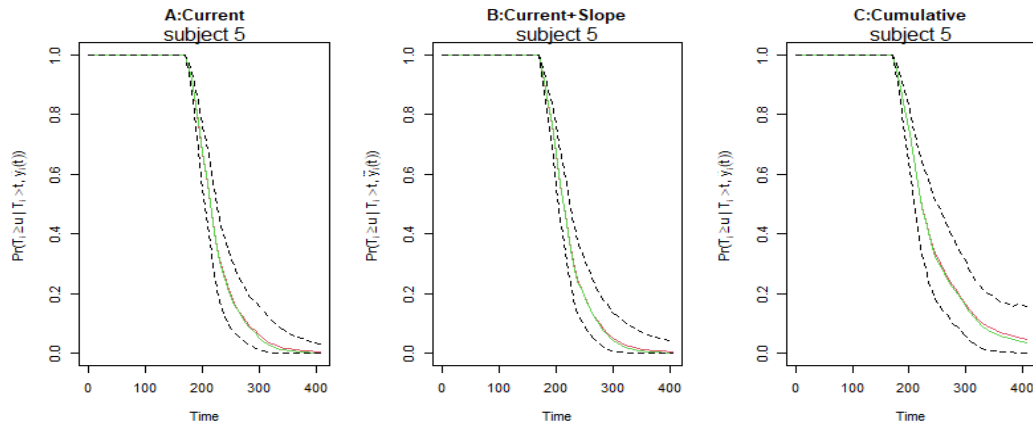
(b) Prediction starting point $t_j = 80\%$

Fig. 5. Reliability function estimates for individual #5 with different association structures at different Prediction starting points. The solid red and green lines depict the mean and median values of the reliability curves, respectively, the black dashed line indicates the 95% confidence interval. (a) at 30% of the lifetime, (b) at 80% of the lifetime.

As expected, the point-by-point interval becomes narrower as the available data increases, The accuracy of the prediction increases as more data becomes available. For comparison, the RUL prediction performance of the joint model with association structure of current value is better than that of the two-stage method when $t_j=30\%$ and 80% , and the prediction results are also shown in Table 5. Again, it validates the effectiveness of our proposed model.

Table 5

Comparison of #5 prediction starting point $t_j=30\%$ and $t_j=80\%$ model RUL prediction results

Prediction starting point t_j	True value of life	Association structure	Bayesian joint model			RUL predicted by two-stage method ^[16]
			Predicted value	2.5%	97.5%	
30%	149	Current	159	144	174	189
		Current+slope	156	143	171	
		Cumulative	167	144	189	
80%	43	Current	45	28	63	53
		Current+slope	44	31	61	
		Cumulative	55	37	74	

5. Conclusion

The work successfully proposed a more accurate RUL prediction method

by jointly analyzing degradation and failure data with association structure identification. In consideration of individual variability and monitoring errors, a linear mixed effects model is applied for the degradation process. To capture the interaction of degradation and failure, a proportional hazards model is utilized to combine the different characteristics of potential degradation processes as built-in covariates and failure time data. The optimal association structure for the joint model is identified based on the prediction performance of validation set. The application for the aircraft engine data sets shows that the joint model with association structure between current degradation value and failure outperforms the two-stage method without association structure in terms of RUL prediction accuracy, especially in the early stage. Moreover, the joint model with current-value association structure or current-slope association structure performs best overall. Therefore, the proposed joint model for RUL prediction with association structure identification may provide reference for other mechatronic equipments.

In this study, the effect of only one monitoring variable on survival outcome was considered, and other monitoring variables were not considered, so there may be many other covariates affecting the risk of failure function. In addition, there may be other association structures such as lagged and random effects to influence the risk of failure function, and these conjectures need to be further analyzed and studied.

Acknowledgment

This research was supported by the Hubei Provincial Department of Education (Grant No. D20221105). This research was also supported by the National Key Research and Development Program of China“Manufacturing Basic Technology and Key Components”key project (Grant No. 2021YFB2011200).

R E F E R E N C E S

- [1]. *Man J, Zhou Q.* Prediction of hard failures with stochastic degradation signals using Wiener process and proportional hazards model. *Computers & Industrial Engineering*, 2018, 125: 480-489.
- [2]. *Ling L.* Residual life prediction for the spindle of retired machine based on fatigue damage mechanism. *UPB Scientific Bulletin, Series D: Mechanical Engineering*, 2018, 80(02): 183.
- [3]. *Ellis B, Heyns P S, Schmidt S.* A hybrid framework for remaining useful life estimation of turbomachine rotor blades. *Mechanical Systems and Signal Processing*, 2022, 170: 108805.

-
- [4]. *Mnoharan P, Dennsion M S, Ganesan V, et al.* Reliability enhancement of steel rolling mill using fault tree analysis. *UPB Scientific Bulletin, Series D: Mechanical Engineering*, 2019, 81(1): 165-178.
 - [5]. *Yin S , Li X , Gao H , et al.* Data-Based Techniques Focused on Modern Industry: An Overview. *IEEE Transactions on Industrial Electronics*, 2015, 62(1):657-667.
 - [6]. *Liu Y Q, Chen Z G, Wang K Y, et al.* Surface wear evolution of traction motor bearings in vibration environment of a locomotive during operation. *Science China Technological Sciences*, 2022, 65(4): 920-931.
 - [7]. *WU Z, LIU Y, LI X, et al.* Intelligent fault diagnosis of robot bearing based on multi-information. *UPB Scientific Bulletin, Series D: Mechanical Engineering*, 2020, 82(3): 145-235.
 - [8]. *Peng C, Chen Y, Chen Q, et al.* A Remaining Useful Life prognosis of turbofan engine using temporal and spatial feature fusion. *Sensors*, 2021, 21(2): 418.
 - [9]. *Saxena A, Goebel K.* C-mapss data set. *NASA Ames Prognostics Data Repository*, 2008.
 - [10]. *Lee M L T, Whitmore G A.* Threshold regression for survival analysis: modeling event times by a stochastic process reaching a boundary. *Statistical Science*, 2006, 21(4): 501-513.
 - [11]. *Liu X, Li J, Al-Khalifa K N, et al.* Condition-based maintenance for continuously monitored degrading systems with multiple failure modes. *IIE transactions*, 2013, 45(4): 422-435.
 - [12]. *Song S, Coit D W, Feng Q.* Reliability analysis of multiple-component series systems subject to hard and soft failures with dependent shock effects. *IIE Transactions*, 2016, 48(8): 720-735.
 - [13]. *PENG WANG, DAVID W. COIT.* Reliability and Degradation Modeling with Random or Uncertain Failure Threshold. //53rd Annual Reliability and Maintainability Symposium (RAMS 2007). 2007:392-397.
 - [14]. *SAE International.* Comprehensive Life Test for 12 V Automotive Storage Batteries. 2013.
 - [15]. *Xiong R, Zhang Y, Wang J, et al.* Lithium-ion battery health prognosis based on a real battery management system used in electric vehicles. *IEEE Transactions on Vehicular Technology*, 2018, 68(5): 4110-4121.
 - [16]. *Hu J, Sun Q, Ye Z S, et al.* Joint modeling of degradation and lifetime data for RUL prediction of deteriorating products. *IEEE Transactions on Industrial Informatics*, 2020, 17(7): 4521-4531.
 - [17]. *Man J, Zhou Q.* Remaining useful life prediction for hard failures using joint model with extended hazard. *Quality and Reliability Engineering International*, 2018, 34(5): 748-758.
 - [18]. *Ge R, Zhai Q, Wang H, et al.* Wiener degradation models with scale-mixture normal distributed measurement errors for RUL prediction. *Mechanical Systems and Signal Processing*, 2022, 173: 109029.
 - [19]. *Alsefri M, Sudell M, García-Fiñana M, et al.* Bayesian joint modelling of longitudinal and time to event data: a methodological review. *BMC medical research methodology*, 2020, 20(1): 1-17.

- [20]. Rizopoulos D, Ghosh P. A Bayesian semiparametric multivariate joint model for multiple longitudinal outcomes and a time-to-event. *Statistics in medicine*, 2011, 30(12): 1366-1380.
- [21]. Han X, Wang Z, Xie M, et al. Remaining useful life prediction and predictive maintenance strategies for multi-state manufacturing systems considering functional dependence. *Reliability Engineering & System Safety*, 2021, 210: 107560.
- [22]. Zhang S, Zhai Q, Shi X, et al. A Wiener Process Model With Dynamic Covariate for Degradation Modeling and Remaining Useful Life Prediction. *IEEE Transactions on Reliability*, 2022.
- [23]. Hennessey V, Novelo LL, Li J, Zhu L, Huang X, Chi E, et al. A Bayesian jointmodel for longitudinal DAS28 scores and competing risk informative dropout in a rheumatoid arthritis clinical trial. 2018.
- [24]. Wang H, Li R, Tsai C L. Tuning parameter selectors for the smoothly clipped absolute deviation method. *Biometrika*, 2007, 94(3): 553-568.
- [25]. Rizopoulos D. The R package JMbayes for fitting joint models for longitudinal and time-to-event data using MCMC. arXiv preprint arXiv:1404.7625, 2014.
- [26]. Li K, Luo S. Dynamic predictions in Bayesian functional joint models for longitudinal and time-to-event data: An application to Alzheimer's disease. *Statistical methods in medical research*, 2019, 28(2): 327-342.
- [27]. Rizopoulos D. Dynamic predictions and prospective accuracy in joint models for longitudinal and time-to-event data. *Biometrics*, 2011, 67(3): 819-829.

Appendix

pseudocode

Algorithm: RUL prediction is based on a joint model with association structure identification between failure time and degradation data.

Input: Training and validation data set, including failure time and degradation data and the whole lifetime. Test data set, degradation data and a part of the whole lifetime.

Output: RUL prediction of the validation set data, RUL prediction error of the validation set data.

Begin:

For $i=1$ to simulation group number N , randomly select sample of 8/9 as training data set $G_{\text{Train}i}$, and remaining is regreded as validation data set $G_{\text{Validation}i}$ (excluding test set).

1) Establish basic framework of joint model

Selecting variables according to AIC/BIC criterion to build degradation model $y_{i,t}(t) = \eta_{i,t}(t) + \varepsilon_{i,t}(t) = X_{i,t}^T(t)\beta + Z_{i,t}^T(t)b_{i,t} + \varepsilon_{i,t}(t)$.

Selecting variables based on LPML/DIC criterion to build a hazard model $h_{i,t}(t) = h_{i0}(t) \exp \left[f \left\{ \eta_{i,t}(t), b_{i,t}, \alpha \right\} \right], t > 0$.

2) Parameter estimation for joint model based on training data set $G_{\text{Train}i}$

Likelihood function $p(y_{i,t}, T_{i,t}, \delta_{i,t} | b_{i,t}, \phi) = p(y_{i,t} | b_{i,t}, \phi) p(T_{i,t}, \delta_{i,t} | b_{i,t}, \phi)$.

Bayesian estimation $\phi = (\beta, D, \sigma, \alpha, \alpha_1, \alpha_2, \alpha_3)$.

For j to starting time number M of RUL prediction for validation data sets

1) Reliability prediction for sample of validation data set $G_{\text{Validation}i}$ from given stating time t_j , $R_{q,i}(t | t_j) = \text{Function2}(t_j, X(t), t \leq t_j)$, JMbayes parameters, $t > t_j$. $\text{Size}(G_{\text{Validation}i}) = K$.

2) RUL prediction for individual λ^{th} of validation data set $G_{\text{Validation}i}$ from given stating time t_j .

3) RUL prediction of t_j at the time of the i th simulation group number in K validation set sample sizes: $\widehat{RUL}_{q,i}(t_j) = \int_0^\infty \widehat{R}_{q,i}(t | t_j) ds$.

4) RUL prediction error for individual λ^{th} of validation data set $G_{\text{Validation}i}$ from given stating time t_j . The i th simulation group number is the prediction error of the remaining useful life at t_j in the sample size of K validation sets: $\text{Error_of_RUL}_{q,i}(t_j) = |\widehat{RUL}_{q,i}(t_j) - RUL_{q,i}(t_j)|$.

5) RUL prediction error for validation data set $G_{\text{Validation}i}$ from given stating time t_j , the average absolute error of the i th analog group number at the t_j moment: $\text{Error_of_RUL}_q(t_j) = \text{MAE}_{ij} = K^{-1} \sum_{\lambda=1}^K |\widehat{RUL}_{q,i}(t_j) - RUL_{q,i}(t_j)|$.

End

End

Estimate RUL prediction error for validation data set for different stating time $t_j, j=1:M$, calculate the median of N average absolute errors at t_j .

$\text{Error_of_RUL}_q(t_j) = \text{median}\{\text{MAE}_{ij}^{(i)}, i = 1: N\}$.

End