FORECASTING REMAINING USEFUL LIFE: INTERPRETABLE DEEP LEARNING APPROACH VIA VARIATIONAL BAYESIAN INFERENCES

2020.10.8 SICHEN LI

GOAL

- Predicting remaining useful life (RUL) for machinery -> preemptive maintenance, prevent failure
- Forecasting via machine learning with "interpretability"
 - the decision logic of the model itself is transparent
 - != post-hoc explainability
- Two types of models
 - Probabilistic lifetime models: easy to interpret, not machine-specific
 - Data-driven ML: black-box, but good predicting power
 - Combine them

PROPOSED STRUCTURED-EFFECT NETWORK

- Decompose the RUL prediction into
 - population-wide baseline: A non-parametric, general lifetime common across all machines
 - machine-specific heterogeneity, including
 - A linear combination of sensor measurement
 - A recurrent component that incorporates historic data



1. COMMON BASELINE: PROBABILISTIC LIFETIME MODELS

- Use past lifetimes of all machines to fit a pre-defined pdf form
- The Weibull distribution, and the log-normal distribution over total lifetime Z: $P_{\text{Weibull}}(Z; a, b) = \frac{a}{b} \left(\frac{Z}{a}\right)^{b-1} e^{-\left(\frac{Z}{a}\right)^{b}}$ and

$$P_{\text{log-normal}}(Z;\,a,\,b) = \begin{cases} \frac{1}{\sqrt{2\pi}bZ}e - \frac{(\log(Z)-a)^2}{2b^2}, \quad Z>0, \end{cases}$$

Conditional expectation, given that the machine had already run for time t

$$\lambda(t) = \mathbb{E}_{Z \sim \text{Weibull}(a,b)}[Z|Z > t] - t \text{ and } \lambda(t) = \mathbb{E}_{Z \sim \text{log-normal}(a,b)}[Z|Z > t] - t$$

Parameter: a, b

2. LINEAR COMPONENT

- input current sensor data: $\beta^T X_t$, or aggregation function over past sensor data: $\beta^T \phi(X_1, ..., X_t)$
- Parameter: beta

Aggregation function	Formula	Interpretation
Max	$\max(X_1,\ldots,X_t)$	Extrema
Min	$\min(X_1,\ldots,X_t)$	Extrema
Mean	$\mu = \frac{1}{2} \sum_{i=1}^{t} X_i$	Average sensor
	t = t = t	measurement
Range	max – min	Variability
Sum	$\sum_{i=1}^{t} X_i$	Total signal
Energy	$\sum_{i=1}^{t} X_i^2$	Total signal with
	$\mathbf{z}_{l=1}$	focus on peaks
Standard deviation	$\sigma = \sqrt{\frac{1}{t} \sum_{i=1}^{t} (X_i - \mu)^2}$	Variability
Skewness	$\frac{1}{t}\sum_{i=1}^{t}\left(\frac{X_{i}-\mu}{\sigma}\right)^{3}$	Symmetry of deviation
Kurtosis	$\frac{1}{2}\sum_{i=1}^{t} \left(\frac{X_i - \mu}{2}\right)^4$	Infrequent extreme
D 1 · 1	$t - t = 1 \left(\sigma \right)$	deviations
Peak-to-peak	$\frac{1}{n_1} \sum_{i=1}^{n_1} \log \max + \frac{1}{n_2} \sum_{i=1}^{n_2} \log \min$	Bandwith
Root mean square	$\sqrt{\frac{1}{t}\sum_{i=1}^{t}X_i^2}$	Total load focus on peaks
Entropy	$-\sum_{i=1}^{t} P(X_i) \log P(X_i)$	Information signal
Arithmetic mean of	$\frac{1}{t}\sum_{i=1}^{t} \operatorname{fft}(X_i)$	Frequency of
power spectral density	$20 \log_{10} \frac{1}{10^{-5}}$	oscillations
Line integral	$\sum_{i=1}^{t-1} X_{i+1} - X_i $	Path length
Kalman filter	$Y_t - b - \sum_{i=1}^p a_i X_{t-i}$	Unexpected
		deviation

3. RNN

Hidden state h1

$$RNN_{\Theta} = f_{NN}([X_t, f_{NN}([X_{t-1}, ... f_{NN}([X_1, 0_n])])])$$

- RNN iterates over the sequence while updating its hidden state ht, which summarizes the already-seen sequence
- The neural network can absorb the variance that cannot be explained by the other components
- Parameter: Theta

MODEL ESTIMATION THROUGH BAYESIAN INFERENCE

• Determine the optimised combined set of unknown parameters $\theta^* = \{a, b, \beta, \Theta\}$ by maximizing the overall likelihood:

$$\theta^* = \operatorname{argmax}_{\theta} P(\theta \mid X) = \frac{P(X \mid \theta) P(\theta)}{P(X)} = \frac{P(X \mid \theta) P(\theta)}{\int P(X \mid \theta) P(\theta) \, \mathrm{d}\theta}.$$

- all parameters with a pre-defined prior distribution
 - Non-parametric component $a \sim \mathcal{N}(a_{empirical,1}), b \sim \mathcal{N}(b_{empirical,1})$
 - Linear component: normal $\beta_i \sim \mathcal{N}(0,10)$; Laplace prior for feature selection
 - Recurrent component: 2-layer LSTM, 100 and 50 neurons; All weights
 ~ Gaussian prior with standard deviation 1

VARIATIONAL BAYES METHOD

- ▶ approximates the true posterior via a variational distribution $Q_{\lambda}(\theta) \approx P(\theta \mid X)$
- Find the optimal lambda*, along with the corresponding distribution Q* that is closest
- Derive a variational lower bound $ELBO(\lambda)$ and do gradient descent to optimise the three components simultaneously $ELBO(\lambda) =: \mathbb{E}_{Q_{\lambda}}[\log P(X, \theta)] - \mathbb{E}_{Q_{\lambda}}[\log Q_{\lambda}(\theta)]$
- Choose λ by $\nabla ELBO(\lambda)$

DATASET

- Turbofan Engine Degradation Simulation dataset
- Predict the RUL (measured in cycles) based on sensor data from 200 aircraft engines -> 100 training, 100 testing
- > 21 sensors

FORECAST RESULT

- Feature-engineering helps
- Traditional ML < Structured-effect network < RNN</p>

Method

MAE

BASELINES WITHOUT SENSOR DATA				
Empirical RU	45.060			
Conditional e	27.794			
Conditional	27.409			
normal)				
TRADITIONAL MA				
Ridge regress	19.193			
Ridge regress	18.382			
engineering)				
Lasso	19.229			
Lasso (with fe	18.853			
Elastic net	19.229			
Elastic net (w	18.245			
engineering)				
Random fore	17.884			
Random for	17.793			
engineer	ring)			
SVR		18.109		
SVR (with fea	21.932			
RECURRENT NEU				
LSTM	11.188			
STRUCTURED-EFF	ECT NEURAL NETWORKS			
Distribution	Linear component			
Weibull	None	15.862		
Weibull	Regularized	17.433		
Weibull	Feature engineering	13.392		
Weibull	Regularized feature	14.989		
	engineering			
log-normal	None	15.061		
log-normal	og-normal Regularized			
log-normal Feature		13.267		
-	engineering			
log-normal	Regularized feature	14.545		
	engineering			

FORECAST DECOMPOSITION



- The distribution-based lifetime component contributes a considerable portion <- overall nature of the RUL (0.175 of variance)
- > The sensor measurements introduce a within-engine and within-time variability (0.408 of variance)
- The recurrent neural network introduces a non-linear black-box component: very small, maybe due to enough predicting power of the current measurement X_t (0.064 of variance)

PAPER



POSTERIOR DISTRIBUTION

Linear components corresponding to 21 sensors

2	U	~ 1	
Sensor	Mean	Standard	Standardized
	estimate	deviation	coefficient
X ₉	- 33.169	0.498	-16.506
X ₁₂	49.721	0.250	12.440
X ₂₁	44.932	0.258	11.601
X ₇	48.154	0.230	11.073
X ₁₁	-24.622	0.357	-8.796
X ₂₀	42.850	0.184	7.880
<i>X</i> ₁₄	-22.540	0.317	-7.155
X_4	-16.183	0.275	-4.447
X ₁₅	-13.934	0.297	-4.145
X_2	-12.446	0.316	-3.931
<i>X</i> ₆	19.678	0.159	3.126
X ₃	-7.851	0.318	-2.494
X ₁₇	- 9.951	0.208	-2.066
<i>X</i> ₈	- 4.121	0.367	-1.511
X ₁₆	-0.273	2.105	-0.574
X ₁₃	-1.424	0.348	-0.495
X ₁₉	0.234	2.028	0.474
<i>X</i> ₁₀	0.126	1.900	0.239
X ₁₈	- 0.095	1.936	-0.184
X_1	-0.070	1.937	-0.135
<i>X</i> ₅	0.001	1.993	0.002

22

21

DERIVATION OF OPTIMISATION PROCEDURE

Appendix A. Derivation of ELBO for structured-effect neural network

Our suggested approach draws upon variational Bayesian methods and approximates the true posterior via a variational distribution $Q_{\lambda}(\theta) \approx P(\theta | X)$. Here $Q_{\lambda}(\theta)$ refers to a family of distributions that is indexed by λ and, hence, our optimization problem translates into finding the optimal λ^* along with the corresponding distribution Q_{λ^*} . The following theorems state the mathematical definition of λ^* and introduce a tractable approximation.

Theorem 1. The optimal λ^* is given by

$$\lambda^* = \operatorname{argmin}_{\lambda} \mathbb{E}_{Q_{\lambda}}[\log Q_{\lambda}(\theta)] - \mathbb{E}_{Q_{\lambda}}[\log P(X, \theta)] + \log P(X).$$
(A1)

Proof. The fit between the variational distribution $Q_{\lambda}(\theta)$ and the posterior distribution $P(\theta \mid X)$ can be measured by the Kullback-Leibler divergence. Hence, we yield

$$\lambda^* = \operatorname{argmin}_{\lambda} KL(Q_{\lambda}(\theta) || P(\theta | X)).$$
(A2)

Inserting the definition of the Kullback-Leibler divergence results into

$$\lambda^* = \operatorname{argmin}_{\lambda} \mathbb{E}_{Q_{\lambda}} \left[\log Q_{\lambda}(\theta) \right] - \mathbb{E}_{Q_{\lambda}} \left[\log P(\theta \mid X) \right] = \operatorname{argmin}_{\lambda} \mathbb{E}_{Q_{\lambda}} \left[\log Q_{\lambda}(\theta) \right] - \mathbb{E}_{Q_{\lambda}} \left[\log P(X, \theta) \right] + \log P(X).$$
(A4)

_

Unfortunately, Eq. (1) is intractable, as it depends on the marginal likelihood of the model, $\log P(X)$. Therefore, the following theorem derives an approximation for the marginal likelihood of the model.

Theorem 2. The marginal likelihood of the modellog P(X) can be approximated by the evidence lower bound, ELBO(λ), i.e.,

$$\log P(X) \ge \mathbb{E}_{Q_{\lambda}}[\log P(X, \theta)] - \mathbb{E}_{Q_{\lambda}}[\log Q_{\lambda}(\theta)] = ELBO(\lambda).$$
(A5)

DERIVATION OF OPTIMISATION PROCEDURE (CONT.)

Proof. Utilizing Jensen's inequality, it holds that

$$\log P(X) = \log \int P(X, \theta) \, d\theta = \log \int P(X, \theta) \frac{Q_{\lambda}(\theta)}{Q_{\lambda}(\theta)} \, d\theta = \log \mathbb{E}_{Q_{\lambda}} \left[\frac{P(X, \theta)}{Q_{\lambda}(\theta)} \right] \ge \mathbb{E}_{Q_{\lambda}} \left[\log \frac{P(X, \theta)}{Q_{\lambda}(\theta)} \right] = \mathbb{E}_{Q_{\lambda}} [\log P(X, \theta)] - \mathbb{E}_{Q_{\lambda}} [\log q(\theta)]. \tag{A7}$$

(A8)

Theorem 3. The optimal λ^* can be approximated by $\lambda^* = \operatorname{argmax}_{\lambda} ELBO(\lambda).$

Proof. From Eqs. (1) and (2), it immediately follows that

$$\lambda^* = \operatorname{argmin}_{\lambda} \log P(X) - ELBO(\lambda).$$
(A9)

As $\log P(X)$ is constant with respect to λ , the value λ^* can be approximated by maximizing ELBO(λ).

In order to optimize *ELBO*(λ), we utilize gradient descent with the gradients defined by

$$\nabla_{\lambda} ELBO(\lambda) = \nabla_{\lambda} \mathbb{E}_{Q_{\lambda}}[\log P(X, \theta)] - \mathbb{E}_{Q_{\lambda}}[\log q(\theta)] = \mathbb{E}_{Q_{\lambda}}[\nabla_{\lambda} \log q(\theta)(\log P(X, \theta) - \log q(\theta))].$$
(A11)

We further utilize Monte Carlo integration to obtain the estimates of the *ELBO*(λ) and the gradient.