

## k-NN을 활용한 터보팬 엔진의 잔여 유효 수명 예측 모델 제안

김정태\*, 서양우, 이승상, 김소정, 김용근  
LIG넥스원 PGM IPS연구소

### A Proposal of Remaining Useful Life Prediction Model for Turbofan Engine based on k-Nearest Neighbor

Jung-Tae Kim\*, Yang-Woo Seo, Seung-Sang Lee, So-Jung Kim, Yong-Geun Kim  
PGM IPS Lab, LIGnex1

**요약** 정비 산업은 사후정비, 예방정비를 거쳐, 상태기반 정비를 중심으로 진행되고 있다. 상태기반 정비는 장비의 상태를 파악하여, 최적 시점에서의 정비를 수행한다. 최적의 정비 시점을 찾기 위해서는 장비의 상태, 즉 잔여 유효 수명을 정확하게 파악하는 것이 중요하다. 이에, 본 논문은 시뮬레이션 데이터(C-MAPSS)를 사용한 터보팬 엔진의 잔여 유효 수명(RUL, Remaining Useful Life) 예측 모델을 제시한다. 모델링을 위해 C-MAPSS(Commercial Modular Aero-Propulsion System Simulation) 데이터를 전처리, 변환, 예측하는 과정을 거쳤다. RUL 임계값 설정, 이동평균 필터 및 표준화를 통해 데이터 전처리를 수행하였고, 주성분 분석(Principal Component Analysis)과 k-NN(k-Nearest Neighbor)을 활용하여 잔여 유효 수명을 예측하였다. 최적의 성능을 도출하기 위해, 5겹 교차검증 기법을 통해 최적의 주성분 개수 및 k-NN의 근접 데이터 개수를 결정하였다. 또한, 사전 예측의 유용성, 사후 예측의 부적합성을 고려한 스코어링 함수(Scoring Function)를 통해 예측 결과를 분석하였다. 마지막으로, 현재까지 제시되어 온 뉴럴 네트워크 기반의 알고리즘과 예측 성능 비교 및 분석을 통해 k-NN 활용 모델의 유용성을 검증하였다.

**Abstract** The maintenance industry is mainly progressing based on condition-based maintenance after corrective maintenance and preventive maintenance. In condition-based maintenance, maintenance is performed at the optimum time based on the condition of equipment. In order to find the optimal maintenance point, it is important to accurately understand the condition of the equipment, especially the remaining useful life. Thus, using simulation data (C-MAPSS), a prediction model is proposed to predict the remaining useful life of a turbofan engine. For the modeling process, a C-MAPSS dataset was preprocessed, transformed, and predicted. Data pre-processing was performed through piecewise RUL, moving average filters, and standardization. The remaining useful life was predicted using principal component analysis and the k-NN method. In order to derive the optimal performance, the number of principal components and the number of neighbor data for the k-NN method were determined through 5-fold cross validation. The validity of the prediction results was analyzed through a scoring function while considering the usefulness of prior prediction and the incompatibility of post prediction. In addition, the usefulness of the RUL prediction model was proven through comparison with the prediction performance of other neural network-based algorithms.

**Keywords** : C-MAPSS dataset, K-nearest neighbor, Principal component analysis, Prognostics health management, Remaining useful life, Turbofan engine

\*Corresponding Author : Jung-Tae Kim(LIGnex1)

email: jungtae.kim2@lignex1.com

Received January 5, 2021

Accepted April 2, 2021

Revised February 17, 2021

Published April 30, 2021

## 1. 서론

과거의 정비기술은 고장 발생 시 수리하는 사후정비(Corrective Maintenance)에 의존하였다. 그러나, 사후정비는 항공기, 발전기와 같은 고가/고안정성 시스템의 경우, 고장 시 발생하는 치명적인 피해로 인해 적용이 어려웠다. 이에 90년대에 이르러서는, 고장 자체를 예방하기 위해 실제 결합 수준과 관계없이 주기적으로 정비를 시행하는 예방정비(Preventive Maintenance)를 시행하였다. 그러나 예방정비는 잦은 중단과 부품 교체로 인해 높은 비용을 발생시켰다[1].

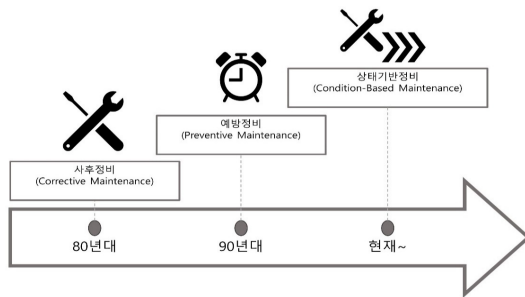


Fig. 1. History of maintenance[2]

현재는 Fig. 1과 같이, 사후정비와 예방정비를 지나 상태기반 정비(Condition-Based Maintenance)를 중심으로 진행되고 있다[2]. 상태기반 정비는 장비의 징후, 상태를 파악하여 고장을 예측하고 능동적인 조치를 취한다. 장비의 전반적인 건전성을 관리하여, 장비 운용을 위한 최적의 시점에 정비가 이루어지고, 불필요한 비용 역시 절감할 수 있다. 따라서 상태기반 정비를 위한 기술 연구가 활발히 이루어지고 있으며, 이 중 대표적인 것이 바로 건전성 예측 및 관리기술(PHM, Prognostics Health Management)이다.

PHM은 기계 설비의 고장을 사전에 예측 진단하여 불시 고장을 방지하고 적절한 유지보수 시기를 예측하는 기술이다[3]. 고장을 사전에 예측 진단하기 때문에 고장 발생 시 치명적 손상이 예상되거나 고안전성을 추구하는 분야에서 다양하게 적용되고 있다.

그 사례로, 차세대 통합 전투기(JSF, Joint Strike Fighter) 개발 컨소시엄은 PHM 프로그램을 수립하여 기술개발을 지속하고 있고, 2003년부터 미국 국방성(DoD, Department of Defense)은 물류확보정책 5000.2에 PHM 시스템을 국방 장비에 의무화하도록 명시하였다[1]. 또한, 미국의 PHM Society는 건전성 예측 및 관리

기술의 발전을 위해 2008년 이래 매년 PHM data competition을 개최하고 있다[4].

관련 연구로, PHM data competition의 터보 엔진의 잔여 유효 수명(이하 RUL, Remaining Useful Life) 예측 주제가 있다. 시뮬레이션 데이터(C-MAPSS Data, Commercial Modular Aero-Propulsion System Simulation) 분석을 통해 잔여 유효 수명을 예측하고, 성능 평가를 위한 점수를 산출하는 연구이다.

본 주제의 연구 사례로, Babu et al.[5]은 컨볼루션 뉴럴 네트워크(CNN, Convolutional Neural Network)를 통해 RUL을 예측하고, MLP(Multi-Layer Perceptron), RVR(Relevance Vector Regression), SVR(Support Vector Regression)과 결과를 비교하였다. Heimes[6]은 시간 순서(time sequence) 데이터를 학습시키는 뉴럴 네트워크인 RNN(Recurrent Neural Network)을 적용하여 RUL을 예측하는 방법을 제시하였다. Yun et al.[7]은 뉴럴 네트워크의 일종인, MLP를 기반으로 RUL을 예측하였다. 또한, Peel, L[8]은 뉴럴 네트워크의 예측 오차를 감소시키기 위해 MLP와 칼만 필터(Kalman filter)를 결합한 모델을 제시하였고, P. Lim et al.[9]는 MLP의 성능을 향상시키기 위해 특징 추출(feature extraction)과 결합한 모델을 제시하였다.

이처럼 해당 주제에 대한 연구는 뉴럴 네트워크 및 그와 결합한 방법론을 활용한 경우가 많다. 그러나 뉴럴 네트워크는 결과값에 대한 투명성을 확보하기 어렵다는 단점이 있다[10]. 연산 과정에 대한 정보를 얻을 수 없으며, 학습 과정에 사용된 모든 매개변수를 찾았다 하더라도, 어떤 과정을 통해 도출되었는지 설명이 어렵다[11].

이에 본 연구에서는 k-NN(k-Nearest Neighbor) 기법을 활용한 RUL 예측 모델을 제시한다. k-NN은 데이터의 군집을 분류하는 분류기(classifier)로, 주변 데이터의 정보를 통해 표적 데이터를 특정 군집으로 분류한다. 데이터 간의 거리를 기반으로 결과를 도출하기 때문에, 연산 과정이 직관적이며 분명하다. 또한, 데이터를 선별하여 활용하기 때문에 이상치나 손실 데이터와 같은 오류 데이터에 민감하지 않다는 장점이 있다. 다만, 본 모델은 기존의 k-NN을 활용하여, 군집 분류가 아닌 표적 데이터의 RUL 값을 예측하였으며, 다른 알고리즘과의 성능 비교를 위해 RUL 예측 정확도에 따른 점수를 산정하였다.

## 2. 본론

### 2.1 분석 데이터(C-MAPSS)

C-MAPSS는 90,000 lb의 추력을 가진 터보 엔진의 시뮬레이션 센서 데이터를 산출하는 MATLAB-Simulink를 말한다. 현재 PCoE(NASA's Prognostics Center of Excellence)에서는 운영조건, 폐쇄 루프 컨트롤러, 고도 및 온도 등 투입(input) 변수에 따른 C-MAPSS 데이터를 Table 1과 같이 4개의 데이터 세트로 구성하여 제공하고 있다[12]. 각 데이터 세트는 고장 유형(Fault Mode), 운영조건(Condition)이 다르며, 학습 및 테스트 데이터로 구분된다.

Table 1. C-MAPSS dataset summary[12]

Dataset	Fault Mode	Conditions	Train Units	Test Units
#1	1	1	100	100
#2	1	6	260	259
#3	2	1	100	100
#4	2	6	249	248

고장 유형은 엔진 구성 요소의 성능 저하 요인들을 말한다. 데이터 세트 1과 2는 한 가지 고장 유형으로 HPC(High Pressure Compressor)의 성능 저하를 갖는다. 반면, 데이터 세트 3과 4는 HPC와 Fan의 성능 저하를 포함한 두 가지 고장 유형을 갖는다. 본 연구에서는 한 가지 고장 유형을 갖는 데이터 세트 1과 2를 대상으로 RUL 예측 모델을 적용한다.

운영조건은 C-MAPSS가 시뮬레이션 데이터를 생성하는데 고려된 운용 환경이다. 비행의 출력에 영향을 주는 데이터 세트 내 운용상의 조건 변수(Operational Setting)를 고려하여 파악한다. 데이터 세트 1은 1개의 운영조건을, 데이터 세트 2는 6개의 운영조건을 갖는다. 따라서, 데이터 세트 2는 운영조건별로 데이터를 분류하여 분석한다. Fig. 2는 데이터 세트 1과 2에 대한 운영조건이다.

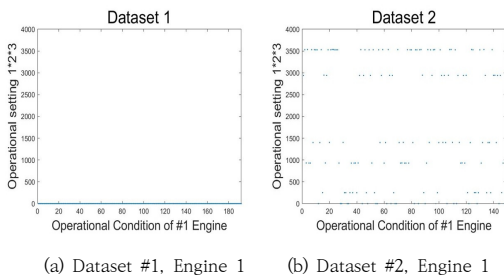


Fig. 2. Operating conditions of the Dataset #1 and #2

Train Units와 Test Units는 학습 및 테스트 데이터를 구성하는 엔진 샘플의 수이다. RUL 예측 모델을 형성할 때, 학습 데이터를 토대로 모델을 구성하고, 테스트 데이터를 적용하여 모델의 성능을 평가한다. 이때, 학습 및 테스트 데이터는 26열로 구성되어 있으며, Table 2와 Fig. 3은 그 데이터 구조이다.

Table 2. Column names of dataset[12]

Column No.	Sensor Value
1	Unit number
2	Time, in cycles
3	Operational setting 1
4	Operational setting 2
5	Operational setting 3
6	Sensor measurement 1
7	Sensor measurement 2
⋮	⋮
26	Sensor measurement 21

Unit	Time	Operational setting			Sensor measurement					
1	2	3	4	5	6	7	8	25	26	
1	1	-7.0000e-04	-4.0000e-04	100	518.6700	641.8200	1.5897e+03	39.0600	23.4190	
1	2	0.0019	-3.0000e-04	100	518.6700	642.1500	1.5910e+03	39	23.4238	
1	3	-0.0043	3.0000e-04	100	518.6700	642.3500	1.5880e+03	38.9500	23.3442	
1	4	7.0000e-04	0	100	518.6700	642.3500	1.5820e+03	38.8800	23.3739	
1	5	-0.0019	-2.0000e-04	100	518.6700	642.3700	1.5829e+03	38.9000	23.4044	
1	6	-0.0043	-1.0000e-04	100	518.6700	642.1000	1.5845e+03	38.9800	23.3669	
1	7	1.0000e-03	1.0000e-04	100	518.6700	642.4800	1.5822e+03	39.1000	23.3774	
1	8	-0.0034	3.0000e-04	100	518.6700	642.5600	1.5830e+03	38.9700	23.3106	
1	9	8.0000e-04	1.0000e-04	100	518.6700	642.1200	1.5910e+03	39.0500	23.4066	
1	10	-0.0033	1.0000e-04	100	518.6700	641.7100	1.5912e+03	38.9500	23.4694	

Fig. 3. Data structure of C-MAPSS dataset

Table 2는 데이터 세트 내 각 열에 대한 정보, Fig. 3은 데이터 세트의 예시이다. 1열의 Unit number는 엔진 번호, 2열의 Cycle은 운전 사이클을 의미한다. 2열의 값을 사용하여 엔진 수명 및 RUL을 계산할 수 있다. 각 엔진의 마지막 사이클 값을 엔진의 수명을 의미하며, RUL은 엔진의 수명에서 현재 사이클 값을 뺀 값이다. 3~5열은 운용조건 파악을 위한 운용상의 조건 변수이다. 운용조건은 3~5열의 조건 변수를 곱한 값이다. 6~26열은 21개의 센서값이다. 데이터 세트 내 각각의 엔진은 21개의 센서를 탑재하고 있으며, Table 3은 해당 센서들에 대한 설명이다.

Table 3. Sensor measurement[12]

Description	Units
Total temperature at fan inlet	° R
Total temperature at LPC outlet	° R
Total temperature at HPC outlet	° R
Total temperature at LPT outlet	° R
Pressure at fan inlet	psia
Total pressure in bypass-duct	psia
Total pressure at HPC outlet	psia
Physical fan speed	rpm
Physical core speed	rpm
Engine pressure ratio	--
Static pressure at HPC outlet	psia
Ratio of fuel flow to Ps30	pps/psi
Corrected fan speed	rpm
Corrected core speed	rpm
Bypass Ratio	--
Burner fuel-air ratio	--
Bleed Enthalpy	--
Demanded fan speed	rpm
Demanded corrected fan speed	rpm
HPT coolant bleed	lbm/s
LPT coolant bleed	lbm/s

## 2.2 모델링

RUL 예측 모델은 데이터 전처리(Preprocessing), 데이터 변환(Transformation), k-최근접 이웃(k-NN) 기법으로 구성하였다. 데이터 전처리는 RUL 임계값 설정, 이동평균 필터, 표준화 과정을 거쳤다. 데이터 변환 과정에서는, 전처리된 데이터를 주성분 분석하여 주성분 맵을 형성하고, k-NN을 활용하여 표적 데이터의 RUL을 산출하였다. 이 과정에서 PHM Society의 스코어링 함수(Scoring function)로 모델 구축을 위한 최적의 파라미터를 검증하고, 모델의 성능을 평가하였다. 전체적인 모델링 절차를 도식화하여 Fig. 4와 같이 제시한다.

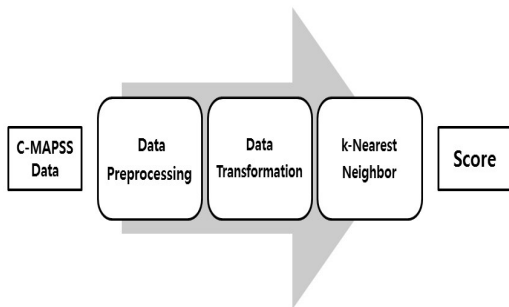


Fig. 4. Modeling Process

### 2.2.1 데이터 전처리

RUL 예측의 목적은 고장을 사전에 예측하여 시기적절한 조치를 취하는 것이다. 따라서, 모델의 관심 대상은 고장 시점에 근접하며, 고장의 전조를 확인할 수 있는 사이클이다. 이를 고려하여 모델을 구성할 때, 학습 데이터의 RUL 임계값을 130 사이클로 설정하였다[6]. 임계값을 초과하는 비관심 데이터의 경우, 수명을 임계값으로 조정하여 분석하였다. Fig. 5는 RUL 임계값 설정을 도식화한 것이다. x축은 Time Cycle, y축은 RUL으로, 사이클이 증가할수록 엔진의 RUL은 감소한다. 엔진의 실제 RUL은 'True RUL'이며, 임계값을 설정한 RUL은 'Piecewise RUL'이다.

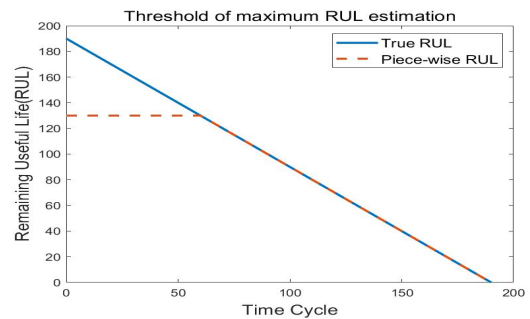


Fig. 5. Threshold of maximum RUL estimation

본 연구에서는 센서 측정값과 운용 사이클의 상관관계를 파악하여 엔진의 수명을 예측하였다. 사이클이 진행되면서 센서 측정값이 상향 또는 하향 추세를 갖는 센서를 분석 대상으로 활용하였다. 따라서, 각 엔진의 21개 센서 측정값 중, 추세를 갖지 않는 7개의 센서를 제외한 14개의 센서를 선정 및 분석하였다. Fig. 6은 센서 측정값의 상향 또는 하향 추세를 나타낸 예시이다.

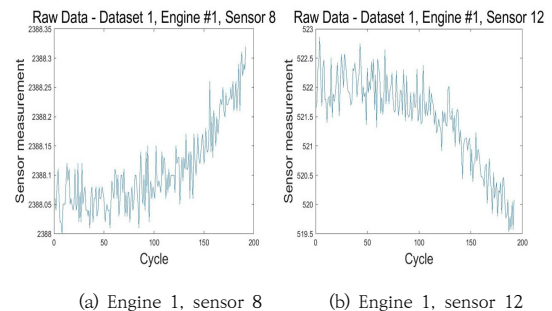
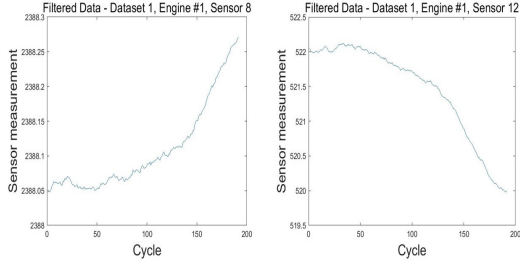


Fig. 6. Raw data of individual engine (Dataset #1)

C-MAPSS 데이터 세트의 센서값 데이터는 Time cycle을 축으로 하는 시계열 데이터이다. 시계열 데이터에 대한 전처리 과정으로, 데이터 이상치와 불규칙 변동 등을 제거하고 데이터의 추세를 반영하기 위해 이동평균 필터를 적용하였다. Fig. 7은 Fig. 6의 Raw data에 이동평균 필터를 적용한 결과이다.



(a) Engine 1, sensor 8

(b) Engine 1, sensor 12

Fig. 7. Filtered data of individual engine (Dataset #1)

각각의 센서값들은 측정 단위 및 척도(scale)가 다르다. 측정 단위 및 척도가 다른 센서값을 같은 척도에서 비교 및 분석하기 위해서는 표준화 과정이 필요하다. 따라서, 엔진 데이터를 통합하고 같은 척도에서 주성분 분석을 하기 위한 표준화 과정을 거쳤다. 표준화 식으로는 Eq. (1)을 적용하였으며, Fig. 8은 필터링된 데이터에 표준화를 적용한 결과이다.

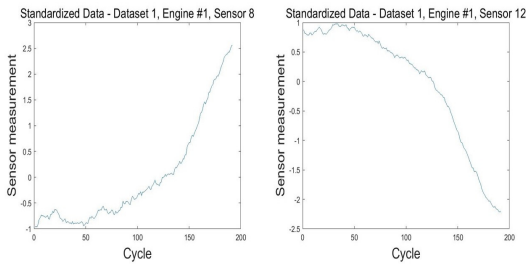
$$x' = \frac{x - \mu}{\sigma} \quad (1)$$

$x'$ : 표준화된 데이터

$x$ : 학습 데이터

$\mu$ : 학습 데이터의 평균

$\sigma$ : 학습 데이터의 표준편차



(a) Engine 1, sensor 8

(b) Engine 1, sensor 12

Fig. 8. Standardized data of individual engine (Dataset #1)

## 2.2.2 데이터 변환

고차원 데이터의 경우 시각화가 어려워, 변수 간의 상관관계를 파악하기 어렵다. 따라서, 주성분 분석을 통해 정보의 손실을 최소화하여 데이터의 차원을 축소하였다.

데이터의 주성분은 다음과 같이 구하였다. 먼저, 14개의 센서를 입력 벡터  $x$ 로 구성하였다. 입력 벡터  $x$ 의 공분산 행렬(Covariance Matrix)을 구하고, 고유값(Eigenvalue)에 따라 고유 벡터(Eigenvector)를 내림차순으로 정렬하였다. 공분산 행렬  $C$ 는 Eq. (2), Eq. (3)을, 고유값과 고유 벡터는 Eq. (4)를 통해 산출하였다.

$$\mu_i = \frac{1}{M} \sum_{i=1}^M x_i \quad (2)$$

$$C = \frac{1}{M-1} \sum_{i=1}^M (x_i - \mu_i)(x_i - \mu_i)^T \quad (3)$$

$M$ : 분석 대상 센서 데이터 개수

$$\det(C - \lambda E) = 0 \quad (4)$$

Fig. 9는 데이터 세트 1의 각  $n$ 번째 주성분( $x$ 축)에 포함된 분산 정도( $y$ 축)를 나타낸 것이다. 즉, 첫 번째 주성분은 데이터 세트 1의 전체 분산 중 약 62%의 분산 정보를 가졌다.

본 연구에서는 주성분을 선택하여 차원이 축소된 데이터들을 Principal Component map(이하 주성분 맵)이라고 정의하였다. 축소된 데이터는 몇 개의 주성분을 사용하는가에 따라 전체 데이터 분산 정보의 일정 부분을 포함한다. 따라서, 검증 과정을 통해 높은 분산 정도를 갖는 순서대로  $m(2 \sim 14)$ 개의 주성분을 선택했을 때 성능을 비교하여 주성분의 개수를 결정하였다.

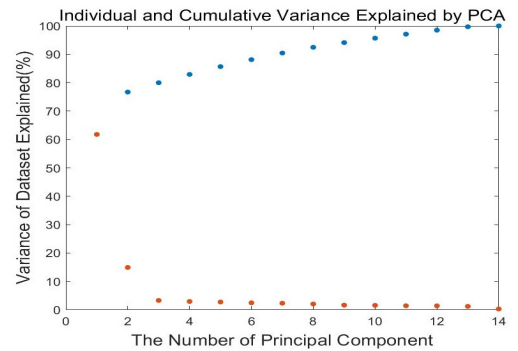


Fig. 9. Individual and cumulative variance explained by Principal Component

Fig. 10은 14차원의 데이터(Dataset#1)를 2개의 주성분을 선택하여 나타낸 예시이다. 청색 점은 전체 엔진 데이터, 적색 점은 1개 엔진 데이터이다. 엔진이 열화되면서 데이터의 위치가 좌측에서 우측으로 이동하는 방향성을 갖는 것을 확인하였다.

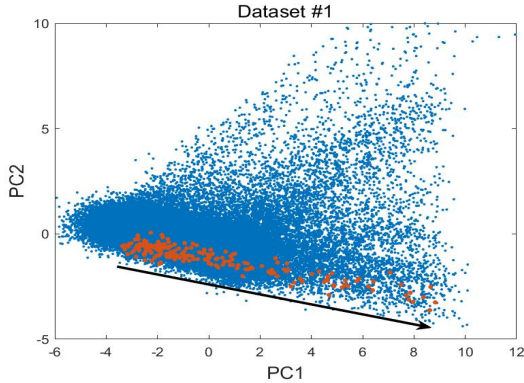


Fig. 10. Example of Principal Component map (Dataset #1, m = 2)

### 2.2.3 k-최근접 이웃(k-Nearest Neighbor)

근접 이웃(NN, Nearest neighbor) 문제는 Beyer, Kevin & Goldstein[13]의 연구에 따르면, “데이터 점의 집합과 미지의 점이 m 차원의 공간에 주어졌을 때, 미지의 점에 가장 가까운 데이터 점을 찾는 것”으로 정의된다. 기존의 k-NN 알고리즘은 주변 데이터 점의 군집을 토대로 분석 대상 데이터를 특정 군집으로 분류하는 것이 목적이다. 따라서, 일반적으로 k-NN 알고리즘의 결과값은 특정 군집이다. 반면, 본 연구에서는 k-NN 알고리즘을 응용하여 테스트 데이터의 RUL을 예측하였다. 테스트 데이터에서 가장 가까운 k개의 주변 데이터를 찾아, 주변 데이터의 RUL 값을 산술평균하여 테스트 데이터의 RUL 값을 산출하였다. 따라서 본 연구의 k-NN 결과값은 특정 군집이 아닌 RUL 값이다.

주성분 맵은 m차원 공간을 형성하며, 각각의 데이터들은 RUL 정보를 포함한다. 거리 측정 방식은 유클리디안 거리(Euclidean’s distance)를 사용하였다. k-NN 알고리즘으로 RUL을 예측하는 과정을 (1)~(5)와 같이 제시한다.

- (1) 학습 세트는 m차원 주성분 맵을 형성한다.
- (2) 예측하고자 하는 데이터(이하 표적 데이터)는 이동 평균 필터, 표준화 과정을 거친다.
- (3) 표적 데이터의 마지막 센서값을 학습 세트의 m개

주축에 사영시킨다.

- (4) 주성분 맵 내, 사영시킨 데이터에서 가장 가까운 k개의 학습 데이터를 찾는다.

- (5) k개 학습 데이터의 RUL 값을 산술평균하고 계산된 RUL을 표적 데이터의 RUL으로 예측한다.

k-NN 알고리즘은 k 값에 따라 알고리즘의 성능이 결정된다. 이에, 5겹 교차 검증 방식으로 최적의 k 값을 검증하였다.

### 2.2.4 검증(Validation)

#### 2.2.4.1 스코어링 함수(Scoring Function)

모델 성능 판단 기준은 PHM Society의 스코어링 함수를 사용하였다[12]. 해당 식은 PHM Society의 알고리즘 성능 평가 기준으로 Eq. (5)와 같다. Table 4는 해당 평가 기준 산정을 위한 3가지 경우이다. 이때,  $RUL_{predicted}$ 은 알고리즘이 예측한 RUL,  $RUL_{real}$ 은 테스트 데이터의 실제 RUL을 의미한다.

Table 4. Cases depending on RUL[12]

$RUL_{real} = RUL_{predicted}$	Case 1
$RUL_{real} > RUL_{predicted}$	Case 2
$RUL_{real} < RUL_{predicted}$	Case 3

$$score = \begin{cases} \sum_{i=1}^n e^{-\left(\frac{d}{13}\right)} - 1 & \text{for } d < 0 \\ \sum_{i=1}^n e^{\left(\frac{d}{10}\right)} - 1 & \text{for } d \geq 0 \end{cases} \quad (5)$$

$$\text{where, } d = RUL_{predicted} - RUL_{real}$$

Case 1은 RUL을 예측하는 이상적인 경우이다. 고장이 발생하는 시점을 정확히 알고 조치를 취할 수 있다. Case 2와 Case 3의 경우, 오차가 같더라도 Case 2가 예측 정비에 더욱 적합하다. Case 2는 고장 예지 시점이 고장 이전이며, Case 3은 고장 이후이다. Case 3의 경우, 사후정비와 다를 바 없어 엔진의 유지보수 작업에 적합하지 않다. 따라서 이를 고려한 스코어링 함수가 Eq. (5)이다. 점수가 작을수록 좋은 예측 성능을 의미한다.

#### 2.2.4.2 5겹 교차 검증(5-fold Validation)

PHM Society의 스코어링 함수를 기반으로, 최적 주성분 개수(m)와 근접 데이터 개수(k) 값 결정을 위해 5겹

교차 검증 기법을 사용하였다. 5겹 교차 검증은 알고리즘의 성능을 평가하기 위해 사용되는 기법이다[14]. 학습 데이터를 5개 그룹으로 나누어 4개 그룹은 학습 세트(training set)로, 나머지 1개 그룹은 검증 세트(validation set)로 설정하였다. 검증 세트를 매 검증마다 다르게 선택하여 총 5번의 검증 결과를 도출하였다. Fig. 11은 세트를 나누는 방식을 도식화한 것이다.

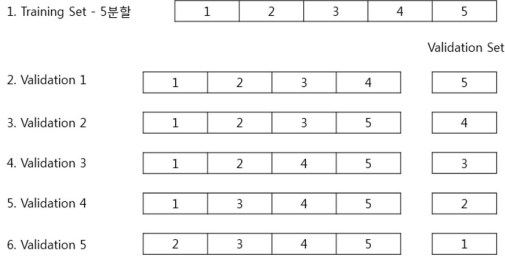


Fig. 11. 5-fold cross validation

검증 대상 파라미터는  $m$ 과  $k$ 이다. 5회 검증 결과를 합산한 점수가 최소인  $m$ 과  $k$  값을 최적값으로 결정하였다. 주성분 개수인  $m$ 은 14개의 센서를 대상으로 2 ~ 14까지 총 13개 값을 사용하였다.  $k$ -NN 알고리즘에서  $k$  값은 범위를 10부터 250까지, 간격은 10으로 하였다. Fig. 12와 Fig. 13은 그 결과이다.

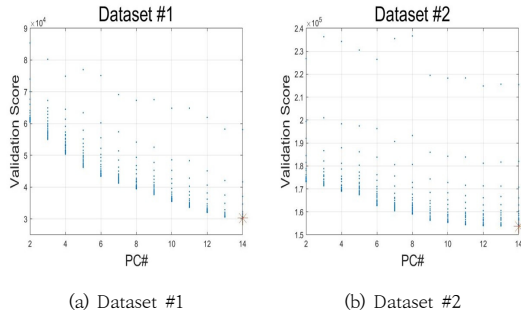


Fig. 12. Validation Score depending on  $m$

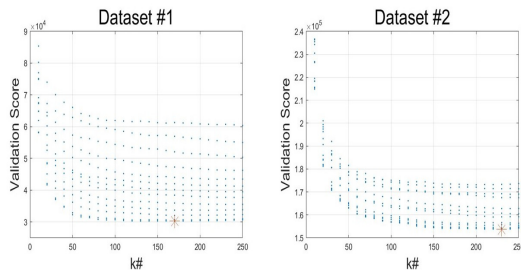


Fig. 13. Validation Score depending on  $k$

Fig. 12는 주성분 개수에 따른 검증 결과, Fig. 13은  $k$  개수에 따른 검증 결과이다. 데이터 세트 1과 2 모두 14개의 주성분을 사용했을 때 가장 좋은 성능을 보였다. 반면,  $k$  개수의 경우, 많을수록 좋은 성능을 보이는 것이 아니라 최적 성능을 갖는  $k$  값이 도출되었다. Table 5는 검증 결과 도출된 최적  $m$ ,  $k$  값과 그 점수이다.

Table 5. Optimal value of  $m$  &  $k$

	Dataset 1	Dataset 2
PC #, $m$	14	14
K #, $k$	170	230
Validation Score	30.273	153.730

### 2.2.5 RUL 예측 모델(RUL Prediction Model)

최적  $m$ ,  $k$  값을 적용하여 RUL 예측 모델을 구성한다. 단, 테스트 세트는 정보가 없을 경우를 대비하여 모델의 모수만으로 RUL을 예측하였다. RUL 임계값 조정 과정은 거치지 않고, 표준화와 주성분 맵은 학습 세트의 모수( $\mu$ ,  $\sigma$ , 공분산)를 사용하였다. Fig. 14는 모델을 도식화한 것이다.

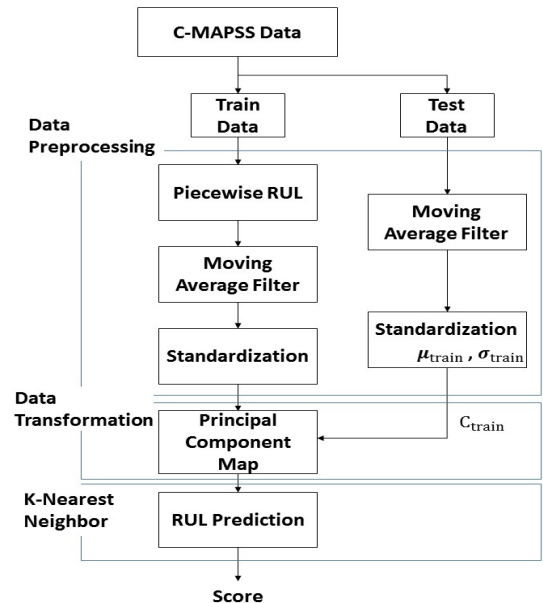


Fig. 14. RUL prediction model

### 2.3 성능 평가

검증을 통해 산출한  $k$ ,  $m$  값을 활용하여 테스트 데이

터에 대한 점수를 산출하였다. 테스트 데이터를 변환하여 점수를 산출하는 과정을 (1)~(5)와 같이 제시한다.

- (1) 테스트 데이터를 이동평균 필터 및 표준화를 통해 전처리한다. 이때, 표준화는 학습 데이터의  $\mu$ 와  $\sigma$ 를 사용하였다.
- (2) 전처리 데이터를 학습 데이터의 주성분으로 변환한다.
- (3) 변환된 데이터를 주성분 맵에 투입한다.
- (4) 각 데이터들 근방의 학습 데이터 k개의 RUL 평균을 계산한다.
- (5) 각 평균을 해당 데이터의 RUL으로 예측하고, Eq. (5)에 대입하여 최종 점수를 산출한다.

Table 6, Fig. 15는 Babu et al.[6]이 비교했던 알고리즘들과 RUL 예측 모델의 점수이다. 결과를 비교해보면, 제시된 5개의 알고리즘 중 가장 좋은 성능을 보이는 것을 확인할 수 있다. 데이터 세트 1의 경우, 811점으로 CNN과 400점 이상의 차이를 보였으며, 데이터 세트 2에서는 10,900점으로 다른 알고리즘들과 2,700점 이상 차이로 월등하게 좋은 성능을 보였다.

Table 6. Score of various methods[6]

Data Set	Dataset #1	Dataset #2
RUL Prediction Model	$8.11 \times 10^2$	$1.09 \times 10^4$
CNN	$1.29 \times 10^3$	$1.36 \times 10^4$
RVR	$1.50 \times 10^3$	$1.74 \times 10^4$
SVR	$1.38 \times 10^3$	$5.90 \times 10^5$
MLP	$1.79 \times 10^4$	$7.80 \times 10^6$

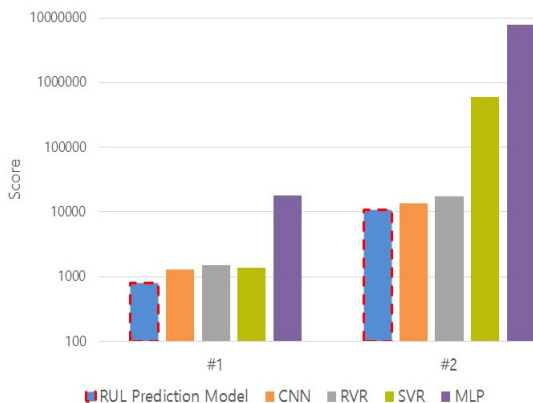


Fig. 15. Score of various methods

### 3. 결론

RUL 예측 모델은 표적 데이터로부터 터보팬 엔진의 RUL을 예측하는 모델이다. RUL 임계값 설정, 이동평균 필터 및 표준화로 데이터를 전처리하고, 주성분 분석을 통해 맵을 형성하였다. 또한, Cross Validation 기법을 활용하여 최적의 주성분 개수(m)와 근접 데이터 개수(k)를 결정하였다. m개 주성분을 선택하여 학습 데이터의 주성분 맵을 형성하고, 테스트 데이터 역시 투입하였다. 마지막으로, 테스트 데이터 근접 k개의 RUL 평균값으로  $RUL_{predicted}$ 를 제시하고 Eq. (5), 스코어링 함수로 성능을 평가하였다. 점수는 낮을수록 실제와 근접한 예측을 의미하며, Table 6에서 확인할 수 있듯이 다른 알고리즘에 비해 모두 좋은 성능을 보였다.

앞서 언급했듯이, C-MAPSS 데이터는 현재까지 뉴럴 네트워크를 중심으로 연구되어왔다. 뉴럴 네트워크는 문제 접근 방식이 일반적이기에 광범위한 문제영역을 다룰 수 있다는 장점이 있다. 반면, 연산 과정에 대한 정보를 얻을 수 없어 결과에 대한 설명력이 부족하기 때문에, 결과값에 대한 투명성을 확보하기 어렵다는 단점이 있다.

이에 반해, 본 연구의 모델은 RUL 예측 과정이 분명하고 직관적이다. 따라서, 센서만 부착되어 있다면 터보팬 엔진뿐 아니라 다른 장비, 장치들에서도 일반적으로 적용할 수 있으며, 인과관계에 대한 파악도 가능하다. 또한, 데이터를 선별하여 활용하기 때문에 이상치나 손실 데이터와 같은 오류 데이터에 민감하지 않다.

더하여, Table 6의 점수를 보면, 데이터 세트 1에서는 400점, 데이터 세트 2에서는 2,700점 이상 좋은 성능을 보임에도 다른 알고리즘에 비해 RUL 예측을 위한 구조가 직관적이고 구현이 쉽다. 또한, 센서기반 장비에 대해 일반적인 적용이 가능하여 범용성을 갖는다. 따라서, 그 자체로도 좋은 성능을 보이며, 일반적인 RUL 예측을 통한 사업성 파악 도구로도 활용될 수 있다. 정밀한 예측을 위한 고비용의 최적화 RUL 예측 알고리즘 개발 전, 대략적인 수명 예측을 통해 대상 장비의 사업성을 판단할 수 있다. 이는 향후 경제적인 측면에서 투자 대비 활용도가 많아 산업현장에서 고려될 수 있는 사항으로 판단된다.

### References

- [1] J. H. Choi, "Introduction of Failure Prediction and Prognostics Health Management", *Journal of the*

*KSME*, Vol.53, No.7, pp.24-34, Jul. 2013.

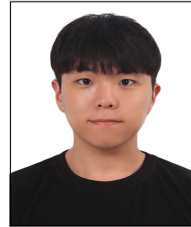
- [2] S. H. Lee, B. D. Youn, "Directions for Industry 4.0, Failure Prediction and Prognostics Health Management", *Journal of the KSME*, Vol.25, No.1, pp.22-28, Feb. 2015.
- [3] B.S. Seo, T. W. Hwang, B. C. Jang, J. H. Song, Y. H. Son, D. K. Lee, B. D. Youn, "Introduction of the 4<sup>th</sup> industrial revolution and success cases through PHM technology", *Journal of the KSME*, Vol.59, No.1, pp.32-37, Jan. 2019.
- [4] X. Jia, B. Huang, J. Feng, H. Cai, J. Lee, "A Review of PHM Data Competitions from 2008 to 2017", *Annual Conference of the PHM Society*, Vol.10, No.1, Sep. 2018.  
DOI: <https://doi.org/10.36001/phmconf.2018.v10i1.462>
- [5] G. S. Babu, P. Zhao, X. Li, "Deep Convolutional Neural Network Based Regression Approach for Estimation of Remaining Useful Life", *Database Systems for Advanced Application. DASFAA*, pp.214-228, Mar. 2016.  
DOI: [https://doi.org/10.1007/978-3-319-32025-0\\_14](https://doi.org/10.1007/978-3-319-32025-0_14)
- [6] Heimes, F., "Recurrent neural networks for remaining useful life estimation", *Intenational Conference on Prognostics and Health Management*, IEEE, Dever, CO, pp.1-6, Oct. 2008.  
DOI: <https://doi.org/10.1109/PHM.2008.4711422>
- [7] Y. Yun, S. Kim, S. H. Cho, J. H. Choi, "Neural Network based Aircraft Engine Health Management using C-MAPSS Data", *Journal of Aerospace System Engineering*, Vol.13, No.6, pp.17-25, 2019.  
DOI: <https://dx.doi.org/10.20910/JASE.2019.13.6.17>
- [8] Peel, L., "Data driven prognostics using a Kalman filter ensemble of neural network models", *Intenational conference on Prognostics and Health Management*, Denver, CO, pp. 1-6, 2008.  
DOI: <https://doi.org/10.1109/phm.2008.4711423>
- [9] P. Lim, C. K. Goh, K. C. Tan, A time-window neural networks based framework for remaining useful life estimation, *International Joint Conference on Neural Networks(IJCNN)*, Vancouver, BC, pp.1746-1753, 2016.  
DOI: <https://doi.org/10.1109/IJCNN.2016.7727410>
- [10] Markus G., Deep learning: a critical appraisal, arXiv:1801.00631, pp.5-14, 2018.
- [11] H. T. Yang, H. Jhang, "Present and future of deep learning", *FUTURE HORIZON*, No.38, pp.8-11, 2018
- [12] A. Saxena, and K. Goebel, PHM08 Challenge Data Set, NASA Ames Prognostics Data Repository NASA Ames Research Center, Moffett Field, CA, 2008.
- [13] K. Beyer, J. Goldstein, R. Ramakrishnan, U. Shaft, "When Is "Nearest Neighbor" meaningful?", *International Conference on Database Theory, ICDT*, pp.217-235, Jan. 1997.  
DOI: [https://doi.org/10.1007/3-540-49257-7\\_15](https://doi.org/10.1007/3-540-49257-7_15)
- [14] T. Wong and P. Yeh, "Reliable Accuracy Estimates from k-fold Cross Validation", *IEEE Transactions on Knowledge and Data Engineering*, Vol.32, No.8,

pp.1586-1594, Aug. 2020.

DOI: <https://doi.org/10.1109/TKDE.2019.2912815>

## 김 정 태(Jung-Tae Kim)

[정회원]



- 2020년 2월 : 한국항공대학교 항공 공기시스템공학과 (공학학사)
- 2020년 1월 ~ 현재 : (주)LIG넥스원 연구원

<관심분야>

건전성 관리, 기계학습

## 서 양 우(Yang-Woo Seo)

[정회원]



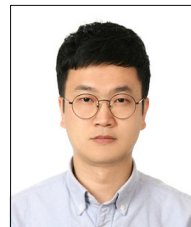
- 1998년 2월 : 홍익대학교 전기공학과 (공학학사)
- 2014년 8월 : 아주대학교 IT융합공학과 (공학석사)
- 2019년 2월 : 아주대학교 시스템공학과 (공학박사)
- 1998년 7월 ~ 현재 : (주)LIG넥스원 수석연구원

<관심분야>

신뢰성, 시스템 엔지니어링, 데이터 분석

## 이 승 상(Seung-Sang Lee)

[정회원]



- 2008년 8월 : 한양대학교 정보경영공학과 (공학학사)
- 2009년 1월 ~ 현재 : (주)LIG넥스원 선임연구원

<관심분야>

기계학습, 신뢰도

---

김 소 정(Kim-So Jung)

[정회원]



- 2018년 2월 : 아주대학교 산업공학과 (공학학사)
- 2020년 2월 : 아주대학교 산업공학과 (공학석사)
- 2020년 1월 ~ 현재 : (주)LIG넥스원 연구원

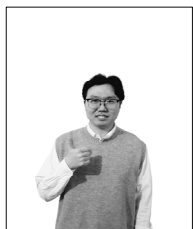
〈관심분야〉

품질, 신뢰성

---

김 용 근(Yong-Geun Kim)

[정회원]



- 2020년 2월 : 한양대학교 산업공학과 (공학학사)
- 2020년 1월 ~ 현재 : (주)LIG넥스원 연구원

〈관심분야〉

신뢰성, 데이터분석