

# 빅데이터 분석을 활용한 초기 정보 기반 화재현장 위험도 예측 모델 개발 연구

## A Study on the Development of a Fire Site Risk Prediction Model based on Initial Information using Big Data Analysis

김도형<sup>1\*</sup> · 조병완<sup>2</sup>Do Hyoung Kim<sup>1\*</sup>, Byung wan Jo<sup>2</sup><sup>1</sup>Ph.D. Candidate, Department of Civil and Environmental Engineering, Hanyang University, Seoul, Republic of Korea<sup>2</sup>Professor, Department of Civil and Environmental Engineering, Hanyang University, Seoul, Republic of Korea

\*Corresponding author: Do Hyoung Kim, dnsura@naver.com

### ABSTRACT

**Purpose:** This study develops a risk prediction model that predicts the risk of a fire site by using initial information such as building information and reporter acquisition information, and supports effective mobilization of fire fighting resources and the establishment of damage minimization strategies for appropriate responses in the early stages of a disaster. **Method:** In order to identify the variables related to the fire damage scale on the fire statistics data, a correlation analysis between variables was performed using a machine learning algorithm to examine predictability, and a learning data set was constructed through preprocessing such as data standardization and discretization. Using this, we tested a plurality of machine learning algorithms, which are evaluated as having high prediction accuracy, and developed a risk prediction model applying the algorithm with the highest accuracy. **Result:** As a result of the machine learning algorithm performance test, the accuracy of the random forest algorithm was the highest, and it was confirmed that the accuracy of the intermediate value was relatively high for the risk class. **Conclusion:** The accuracy of the prediction model was limited due to the bias of the damage scale data in the fire statistics, and data refinement by matching data and supplementing the missing values was necessary to improve the predictive model performance.

**Keywords:** Fire Site Risk, Prediction Model, Big Data Analysis, Machine Learning Algorithm, Random Forest

### 요약

**연구목적:** 본 연구는 화재발생 건축물 정보, 신고자 취득 정보 등 초기 정보를 활용하여 화재현장의 위험도를 예측하여, 재난 발생 초기에 효과적인 소방자원 동원 및 적절한 대응을 위한 피해최소화 전략 수립을 지원하는 위험도 예측 모델을 개발하고자 한다. **연구방법:** 화재 통계 데이터 상에서 화재의 피해규모와 관련된 변수 규명을 위해 머신러닝 알고리즘을 이용한 변수간 상관성 분석을 실시하여 예측 가능성을 검토하고, 데이터 표준화 및 이산화 등의 전처리를 통해 학습 데이터 셋을 구축하였다. 이를 활용하여 예측 정확도가 높은 것으로 평가 받고 있는 복수의 머신러닝 알고리즘을 테스트하여 가장 정확도가 높은 알고리즘을 적용한 위험도 예측 모델을 개발하였다. **연구결과:** 머신러닝 알고리즘 성능 테스트 결과 랜덤포레스트 알고리즘의 정확도가 가장 높게 나왔으며, 위험도 등급에 대해서는 중간치에 대한 정확성이 상대적으로 높은 것으로 확인되었다. **결론:** 화재 통계상 피해규모 데이터의 편향성에 의해 예측모델 정확도가 제한적으로 나타났으며, 예측 모델 성능 개선을 위해 데이터 정합성 및 결손치 보완 등을 통한 데이터 정제가 필요하다.

**핵심용어:** 화재현장 위험도, 예측 모델, 빅데이터 분석, 머신러닝 알고리즘, 랜덤포레스트

Received | 15 March, 2021

Revised | 22 April, 2021

Accepted | 30 April, 2021

 OPEN ACCESS


This is an Open-Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/4.0>) which permits unrestricted noncommercial use, distribution, and reproduction in any medium, provided the original work is properly cited.

© Society of Disaster Information All rights reserved.

## 서론

### 연구배경 및 필요성

국가적 차원의 관심과 지원, 그리고 소방조직의 노력에 의해 국내의 화재대응 능력은 많은 발전을 이루어 왔으나, 시설 노후화, 사회기반시설 밀집 증가 등에 따른 재난사고 취약성은 날로 증가하고 있는 실정이다. 또한 고밀도 도심화, 건축물의 초고층화, 지하공간 증가 등 생활환경 변화에 따른 신종재난의 발생 가능성이 증가하고 있으며, 기술화·산업화 등으로 인한 유해화학물질 유출사고 증가와 같은 특수재난이 증가하고 있다. 가연성 가스 등 고위험물 산재로 인한 현장대원의 재난대응 여건이 악화되고, 초고층 빌딩, 대규모 다중이용시설 등의 증가로 기존 소방기술 적용 어려움이 증가하고 있으며, 시가지 밀집화, 다용도 시설 증가 등으로 인한 재난현장 정보의 부족 문제가 발생하고 있다.

### 연구목적

이런 재난여건 다변화로 인하여 일선 재난현장에서 화재진압, 긴급구조업무를 담당하고 있는 소방조직의 현장 대응역량의 중요성이 늘어나고 있으나, 현장 위험정보의 부족으로 기존의 경험에 의존한 현장대응체계의 한계점을 노출시키고 있다.

화재를 포함한 재난으로 인한 피해를 최소화하기 위해서는 초기 대응이 무엇보다 중요하며, 적절한 대응 방법 및 필요한 소방력과 재난자원에 대한 결정, 동원 등이 신속하게 이뤄졌을 때 그 효과가 제대로 발휘될 수 있다. 그러나 재난현장에는 피해 규모에 영향을 미치거나 확대시킬 수 있는 다양한 상황변수가 존재하고 있어 모든 요소들을 고려하여 사전에 재난현장의 피해규모를 예측한다는 것은 매우 어려운 일이다. 따라서 재난 발생 초기에 효과적인 피해최소화 전략을 수립하고, 최적의 소방력 투입을 위해 초기정보를 기반으로 해당 화재 현장에서 발생한 화재규모를 어느 수준까지 예측하는 것이 가능한지에 대한 검토가 우선적으로 선행되어야 한다. 또한 이를 소방의 현장대응활동에 활용 가능한 수준으로 정확성을 향상시키고, 실제 현장에서 사용하는 정보시스템과 연계하여 활용하는 방안에 대한 검토가 이루어져야 한다.

### 연구방법

본 논문에서는 초기 정보만을 가지고 화재현장의 위험도 예측을 위한 모델을 개발하기 위해, 최근 11년간 S시에서 발생한 화재관련 통계 자료를 조사하여 빅데이터 활용을 위한 기초 및 품질 분석을 실시하였다. 이를 토대로 통계 데이터 상의 각 변수 간 선형 상관성 분석, 머신러닝 알고리즘을 이용한 빅데이터 분석을 실시하였으며, 데이터 정합성 및 결손치에 대한 보완을 위해 표준화 및 이산화 등을 통한 학습 데이터를 구축하였다. 최종적으로는 화재현장 위험도를 정의하여 등급화하고, 알고리즘 성능 테스트를 통해 화재통계 데이터를 활용한 위험도 예측에 가장 적합한 알고리즘을 선정하여 예측 모델을 개발하였다.

## 화재통계 데이터 기초 및 품질 분석

### 분석 개요

과거 11년간의 화재통계 데이터를 이용해서 화재현장 위험도를 예측하는 모델을 개발하기 위해 화재통계 데이터의 기초 및 품질 분석을 수행하였다. 화재통계 데이터는 2009년~2019년의 11년간 S시에서 발생한 화재 가운데 64,647건의 화재에

대한 자료를 분석하였다. 화재통계 데이터의 관리정보 항목으로는 인명피해, 소실면적, 재산피해 등의 화재피해 결과에 관련된 항목과 날씨, 온도, 습도, 풍속 등 화재 규모에 영향을 주는 자연원인과 출동소요시간, 소방서와의 거리, 안전센터와의 거리 등의 사회적 원인 등 총 216개의 관련 정보 항목으로 구성되어 있다.

## 분석 결과

화재통계 상의 216개 항목을 모두 분석하였으며, 본 논문에서 소개할 주요 항목에 대한 분석 결과는 Fig. 1과 같다.

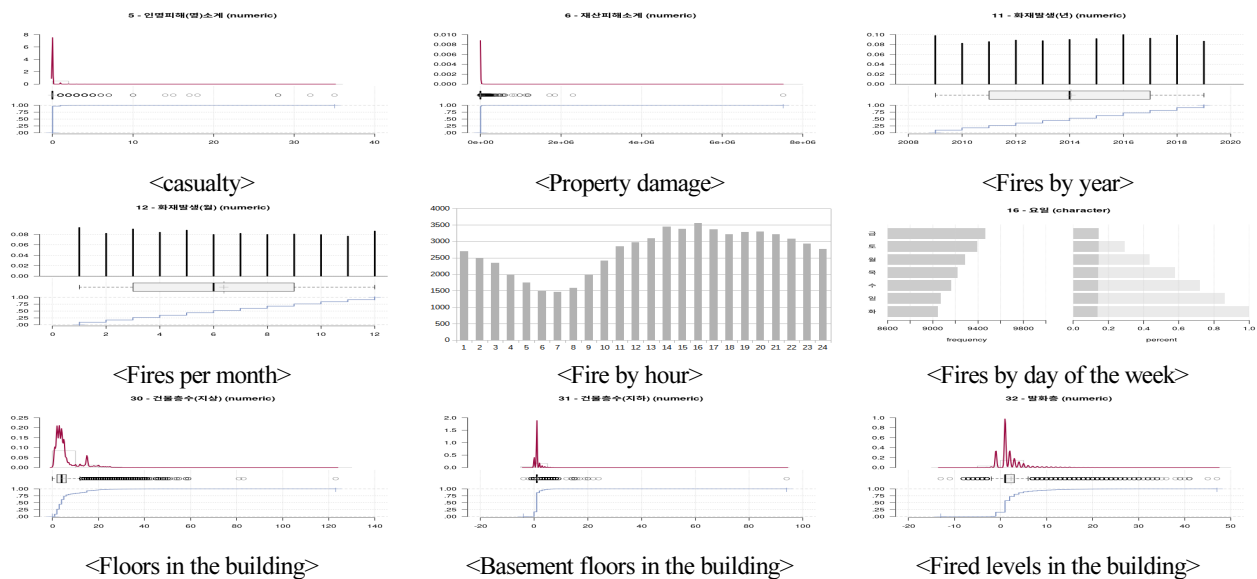


Fig. 1. Results of basic analysis of fire statistics data

먼저, 인명피해 항목을 보면 96.8%인 62,592 건이 인명피해가 0으로 나타나며, 가장 높은 값은 35명, 평균은 0.04명으로 나타났다. 인명피해 분포차트에서 보면 거의 모든 정보가 0 근처에 몰려 있고, 그 이외의 정보는 이상치(outlier)로 인식되고 있다. 인명피해가 거의 0에 몰려 있는 점에서 인명피해를 설명변수로 하는 화재 위험도 예측모델 생성에서 설명변수의 편차(bios) 때문에 예측모델의 성능이 제한적일 것으로 판단되며, 모델 성능 향상을 위해서 다른 설명변수와의 조합을 통한 변수 합성 방법이 필요할 것으로 판단된다. 재산피해의 경우 5.7%가 0이며 평균 2,421.83원, 가장 큰 금액이 7,512,359.0원으로, 대부분의 데이터가 0에 근접해서 대부분은 경미한 재산피해로 화재가 진압되나 일부 케이스에서는 큰 피해가 발생했음을 알 수 있다. 화재발생 건수의 경우 조사기간 동안 크게 변동없이 매년 6천여건 정도가 발생하고 있으며, 월별 화재발생 건수를 보면 겨울~봄 기간인 12월~5월까지가 나머지 기간보다 상대적으로 발생건수가 큰 것을 알 수 있으나, 차이가 크지는 않음을 알 수 있다. 화재발생 시간대 분포는 16시에 피크를 이루다가 점차 감소하다가 07시에 최저점을 지나서 다시 상승하는 패턴을 나타냈으며, 요일별 화재발생 분포를 보면 금, 토, 월요일 순으로 나타나며, 일요일과 화요일이 가장 낮게 나타나고 있다. 화재가 발생한 건물의 평균은 5.8층이며, 최대 123층이며, 건물층수 분포 차트에서 보면 15층 부근에서 빈도가 높아지는 것은 아파트의 층수가 15층 전후가 많기 때문인 것으로 판단된다. 건물의 층수는 건물의 규모와 연관이 되며, 이는 화재 진압

에 있어서 중요한 요소일 것으로 판단되는데, 결손치 비율이 23.5%로 높음으로 인해서 예측모델에 적용하더라도 그 효과가 제한적일 것으로 판단된다.

## 화재통계 데이터 빅데이터 분석

화재현장의 위험도 분석을 위해서 화재통계 데이터에 대해서 화재의 원인과 관련된 변수와 화재의 결과와 관련된 변수를 분별하고 이들 간의 관계 규명을 통해서 화재의 피해 예측가능성에 대해서 검토하였다.

### 각 변수 간 선형 상관성 분석

분석에는 우선 각 변수 간의 상관분석을 바탕으로 화재의 원인과 결과에 관련된 변수들 간의 관계를 고찰하며, 선형적인 분석을 수행하였다. 본 상관분석은 화재 규모의 원인이 되는 항목과 화재로 인한 결과 간의 상관관계를 규명하기 위한 것이며, 상관분석을 통해서 연관성이 높은 속성의 관계성을 규명하고자 하였다. 수집된 자료는 다양한 형태의 데이터이며, 그 가

**Table 1.** Fire statistics data correlation analysis result

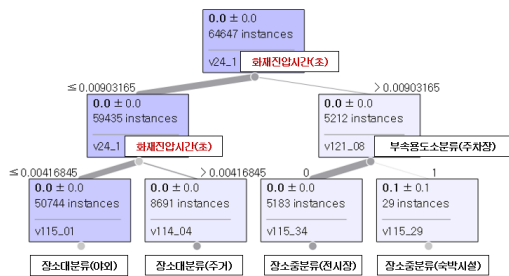
변수1(항목)	변수2(항목)	상관계수	해설
건물층수(지상)	층수(지상)	1	동일의미
건물층수(지하)	층수(지하)	1	동일의미
발화층	발화층수지상(층)	1	동일의미
화재피해경감액	부동산경감	0.998	결과 - 결과
부상	인명피해(명)소계	0.985	결과 - 결과
소실면적(m <sup>2</sup> )	일반직동원명수	0.965	결과 - 결과
전체소실건축구조물(동)	부분소실건축구조물(동)	0.95	전체 - 부분
재산피해소계	동산	0.939	전체 - 부분
재산피해소계	부동산	0.937	전체 - 부분
동원장비소계	펌프물탱크	0.903	전체 - 부분
소실면적(m <sup>2</sup> )	경찰동원명수	0.888	결과 - 결과
전체소실차량등(대)	부분소차량등(대)	0.801	전체 - 부분
동원인력소계	소방동원명수	0.768	전체 - 부분
부동산	동산	0.76	전체 - 부분
소방동원명수	동원장비소계	0.753	결과 - 결과
동원장비소계	구조	0.741	전체 - 부분
층수(지상)	발화층수지상(층)	0.687	무의미
건물층수(지상)	발화층수지상(층)	0.687	무의미
전체소실차량등(대)	전소차량등(대)	0.677	전체 - 부분
동원장비소계	구급	0.674	전체 - 부분
소실면적(m <sup>2</sup> )	동원인력소계	0.672	결과 - 결과
동원인력소계	경찰동원명수	0.669	전체 - 부분
소방동원명수	펌프물탱크	0.669	전체 - 부분
전기가스유관기관동원명수	헬기	0.653	결과 - 결과

운데는 결손치도 상당하여, 수치형 변수에 대해서만 상관분석을 수행한 후 0.6 이상의 상관계수를 가진 데이터 페어에 대한 의미 분석을 수행하였다.

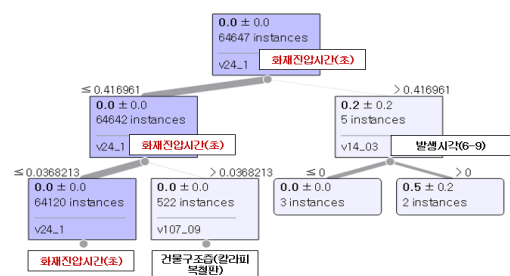
분석결과, 초기에 목적인 바와 같이 화재나 화재로 인한 피해 규모와 그 원인 간의 관계에 있는 변수 간의 상관계수는 무의미한 수준이어서, 의미 있는 결과 도출에 실패하였다. Table 1과 같이 높은 상관계수를 나타내지만 두 변수가 실제 같은 데이터를 나타내는 동일한 의미와 발견하고자 하는 화재 원인-결과 관계가 아닌 결과-결과 관계이거나, 전체와 부분을 나타내는 전체-부분관계 및 실질적으로 전혀 의미가 없는 무의미한 관계 등이 주로 나타났다.

### 머신러닝 알고리즘을 이용한 상관성 분석

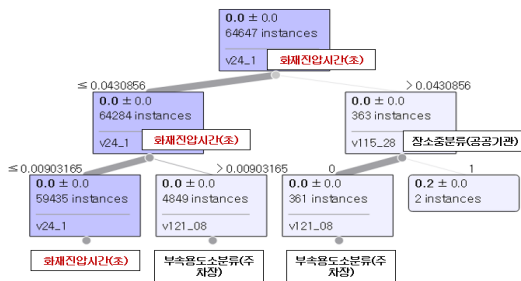
앞 절에서 선형 상관분석을 통해서 화재통계 데이터에서 화재의 피해를 직접적으로 설명할 수 있는 상관성이 높은 변수가 존재하지 않음을 확인하였다. 이에 목적변수를 설명변수로 사용하는 의사결정나무 알고리즘을 이용하여, 화재의 피해에 관련된 변수에 대한 의사결정나무 모델을 Fig. 2와 같이 생성하여 검토하였다. 피해 관련 변수인 인명피해 데이터에 대한 의사결정나무 모델에서 화재진압시간(초) 변수가 가장 중요한 변수임을 알 수 있으나, 이 변수는 화재가 종료된 후에 결정되는 것이므로 화재 피해를 설명할 수 있으나 예측에는 사용할 수 없음을 확인하였다. 다음으로 재산피해 데이터에 대한 의사결정나무 모델링 결과에서도 화재진압시간이 가장 결정적인 역할을 하는 것으로 나타나, 인명 및 재산피해 의사결정나무 모델과 동일하게 예측에는 활용하기 어려운 문제가 있음을 확인하였다. 인명, 재산, 이재의 모든 피해관련 변수에 대해서 각 목적변수를 설명하는 설명변수에 화재진압시간이 공통적으로 나타났다으며, 이 화재진압시간을 결정하는 설명변수를 확인하기 위해서 설명변수들만으로 의사결정나무 모델을 생성하였다. 화



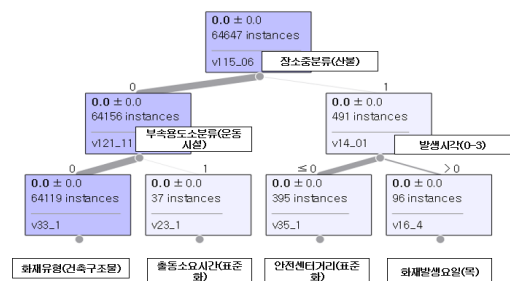
<Decision tree model of casualty variables>



<Decision tree model of property damage variables>



<Decision tree model of variables of casualty and property damage>



<Decision tree model of fire suppression time variable>

Fig. 2. Decision tree model analysis results for fire damage variables



외하였다. 둘째, 변수 중에서 건물화재와 직접적인 관련이 없는 보험, 임야 관련 변수를 제외하였다. 셋째, 화재발생 시각의 영향을 모델에 반영하기 위해서 화재발생 시각을 3시간 단위로 분류한 후에 화재발생 시각에는 1을, 전후 시간대는 0.5를 부여하여 발생시간대의 영향력이 모델에 반영되도록 하였다. 넷째, 요일 변수를 각 요일 별의 화재발생 가능성을 고려해서 요일(월), 요일(화)...요일(일)의 7개의 변수로 분류하고 화재 케이스별로 발생 요일은 1, 이외는 0으로 데이터를 변환하였다. 다섯째, 발화층수의 경우 원 데이터가 지상과 지하가 같은 항목으로 잡혀 있는데, 이를 지하 10층, 지하 5층, 지하 1층, 지상 1층, 지상 5층, 지상 10층, 지상 15층 단위로 변수를 이산화하였으며, 각 화재케이스 별로 발화층이 포함된 변수는 1, 이외는 0으로 이산화하였다. 이상의 데이터 변환을 통해서 Fig. 3과 같이 설명변수 30개에서 234개의 변환변수를 도출하였다.

### 화재현장 위험도 정의 및 등급화

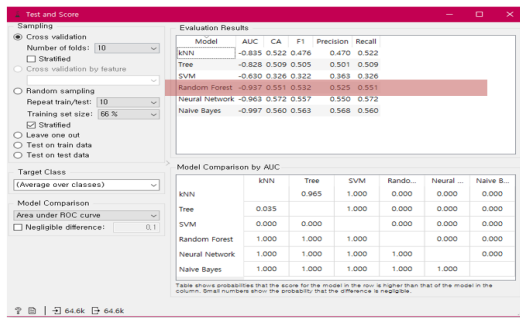
화재통계 데이터에서 피해결과로 정의된 4가지 변수(인명피해, 재산피해, 이재민수, 이재세대수)를 조합해서 Table 1의 산정식과 같은 가중치로 화재 위험지수를 정의해서 하나의 연속값 변수를 도출하였다. 화재 위험지수는 연속적인 수치로서 이를 위험도에 따라서 균등하게 분배하기 위해서 4분위수의 최소값 이하는 0으로 하고, 각 분위별로 값을 1씩 올려서 각 변수를 0 ~ 4의 범위로 변환하는 방식으로 경계를 구분하여 Table 2와 같이 화재위험도를 등급화하였다.

Table 2. Fire site risk definition and rating

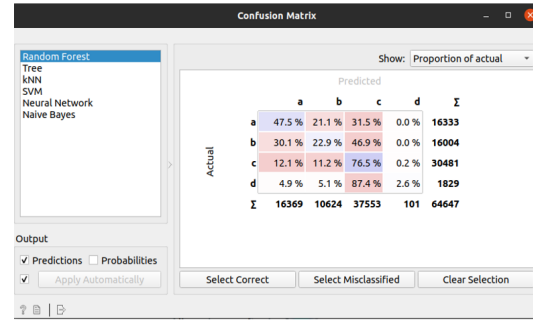
화재 위험지수 산정식	4분위수 개념	화재위험도 등급	데이터 수
$0.4 \times (\text{사망자} \times 0.8 + \text{부상자} \times 0.2) + 0.4 \times \text{재산피해} + 0.2 \times \text{이재피해}$		A(위험-하)	34475
		B(위험-중)	18804
		C(위험-상)	11368

### 머신러닝 알고리즘 성능 테스트

예측모델은 데이터와 알고리즘의 적합성에 의해서 성능이 결정되므로, 데이터 변환을 완료한 학습데이터에 복수의 알고리즘을 적용하여 적합한 알고리즘을 발견하고자 하였다. 테스트 대상 알고리즘으로는 kNN(k-nearest neighbors), decision tree algorithm, SVM(Support Vector Machine), Random Forest, Neural Network, Naive Bayes의 6개의 알고리즘을 선정하였다. 테스트 결과 Fig. 4와 같이 Random Forest, Neural Network, Naive Bayes algorithm이 정확도가 높게 나왔으며, 본 연구에서는 이 3개의 알고리즘 중 확장성 및 정확성 등을 고려해서 Random Forest 알고리즘을 사용하여 모델을 개발하였다. Random Forest 알고리즘을 이용하여 학습 데이터에 개별 등급 별 정확도를 산출하여 검증한 결과, 전체 데이터에 대해 위험도 등급별로 A등급은 47%, B등급은 22.9%, C등급은 76.5%의 정확도를 보였으며, B등급이 상대적으로 정확도가 낮게 나타났다. 이는 인명 및 재산 피해 데이터의 편향성에 의해서 화재위험도의 정확도가 제한적으로 나타난 것으로, 데이터 정합성 및 결손치에 대한 보완을 통해서 데이터 정제 및 모델 개선이 필요할 것으로 판단되었다.



<Algorithm performance comparison result>



<Prediction accuracy verification(random forest)>

Fig. 4. Machine learning algorithm performance test results

## 결론

본 논문에서는 과거 11년간의 S시 화재통계 데이터를 분석하여 화재 원인 및 피해결과와 관련된 변수를 분별하고, 머신러닝 알고리즘을 이용한 상관성 분석 등 빅데이터 분석 결과를 기반으로 한 화재 위험도 예측의 가능성을 검증하였다. 또한 예측 정확도 향상을 위해 화재통계 데이터의 표준화 및 이산화를 통해 학습 데이터 셋을 구축하였으며, 이를 활용하여 화재 위험도 예측에 적합한 알고리즘을 선정하고 예측 모델을 개발하였다. 개발된 예측모델의 정확도를 검증한 결과, 데이터의 편향성에 의해 화재위험도의 정확도가 제한적으로 나타난 일부 등급을 제외하고는 유의미한 결과를 확인할 수 있었다. 예측모델의 정확성 향상을 위해서는 데이터 정합성 및 결손치에 대한 보완을 통해서 데이터를 정제하고, 화재 위험도 예측과 관련있는 신규 변수 발굴 등을 통한 모델 개선이 필요할 것으로 판단되었다. 본 연구를 통해 개발된 예측모델은 화재 발생에 따라 수집되는 정보에 맞춰 예측결과가 지속적으로 업데이트되며 현장상황 변화를 실시간으로 반영할 수 있도록 설계되어 실시간 현장 데이터가 올라오는 소방조직 정보시스템과 연동하여 활용이 가능하며, 통계 기반 예측시스템의 한계 극복을 위해 현장 경험에서 우리나라의 전문가의 식견을 룰셋 기반으로 알고리즘에 반영하는 전문가 시스템 적용에 대한 추가적인 연구가 필요할 것으로 판단된다.

## References

- [1] Chai, S.S., Jang, S.Y., Suh, D. (2018). "Design and implementation of big data analytics framework for disaster risk assessment." Journal of Digital Contents Society, Vol.19, No. 4, pp. 771-777.
- [2] Choi, J.H., Lee, S.W., Hong, W.H. (2013). "A development of fire risk map and risk assessment model for urban residential areas by raking fire causes." Journal of the Architectural Institute of Korea Planning & Design, Vol. 29, No. 1, pp. 271-278.
- [3] Choi, S.H. (2010). "Natural disaster damage cost prediction model based on neural network and genetic algorithm." Proceedings of the Korean Information Science Society Conference, Seoul, pp. 380-384.
- [4] Jeong, M.G., Lee, S.H., Kim, C.S. (2020). "A study on the safety index service model by disaster sector using big data analysis." Journal of the Society of Disaster Information, Vol. 16, No, 4, pp. 682-690.
- [5] Ko, K.S., Hwang, D.H., Park, S.J., Moon, G.G. (2018). "Electrical fire prediction model study using machine learning." The Journal of Korea Institute of Information, Electronics, and Communication Technology, Vol. 11, No.



6, pp. 703-710.

- [6] Kwon, Y.J., Kim, D.E. (2009). "A study on the development of evaluation methods for fire risk analysis of high-rise building." *Proceedings of the Korea Institute of Fire Science and Engineering Conference, Busan*, pp. 270-275.
- [7] Lee, C.Y., Kim, T.H., Cha, S.Y. (2011). "A study of the extraction algorithm of the disaster sign data from web." *Journal of The Korean Society of Disaster Information*, Vol. 7, No. 2, pp. 139-149.
- [8] Lim, H.S., Lee, K.M., Cho, J.W., Lee, S.K., Min, S.H., Min, J.K. (2019). "Analysis of the relationship between fire factors and influential factors using SPSS." *Journal of the Korean Society of Hazard Mitigation*, Vol. 19, No. 5, pp. 103-112.
- [9] Park, E., Min, S. (2019). "Standardization of fire factor for big data." *Journal of the Korean Society of Hazard Mitigation*, Vol. 19, No. 4, pp. 143-149.
- [10] Ryu, J.W., Kwon, S.P. (2015). "Fire risk assessment based on weather information using data mining." *Fire Science and Engineering*, Vol. 29, No. 5, pp. 88-95.
- [11] Shin, J.D., Jeong, S.H., Kim, M.S., Kim, H.J. (2012). "Analysis of fire risk with building use type using statistical data." *Journal of The Korean Society of Hazard Mitigation*, Vol. 12, No. 4, pp. 107-114.
- [12] Shin, Y.C., Koo, I.H., Hayashi, Y., Ohmiya, Y., Kwon, Y.J. (2011). "A study on the risk assessment using simulation and case study of urban fire-focusing on market." *Fire Science and Engineering*, Vol. 25, No. 6, pp. 1-7.