

http://dx.doi.org/10.17703/JCCT.2024.10.5.725

JCCT 2024-9-86

화학물질 사고 현황 및 사례 데이터를 이용한 인공지능 사고 원인 예측 모델에 관한 연구

A Study on Artificial Intelligence Models for Predicting the Causes of Chemical Accidents Using Chemical Accident Status and Case Data

이경현*, 백락준**, 정혜성***, 김우수****, 최희정*****

KyungHyun Lee*, RackJune Baek**, Hyeseong Jung***,
WooSu Kim****, HeeJeong Choi*****

요약 본 연구는 환경부 산하 화학물질안전원에서 제공하는 2014년 1월부터 2024년 1월까지의 화학물질 사고 현황 및 사례 데이터 865건을 활용하여 인공지능 기반 사고 원인 예측 모델을 개발하는 것을 목표로 한다. 본 연구에서는 6개의 인공지능 모델을 사용해 데이터를 학습시키고, 평가지표로는 정확도(Accuracy), 정밀도(Precision), 재현율(Recall), F1 스코어(F1 Score)를 비교 분석하였다. 2020년부터 2024년까지의 화학물질 사고 사례 데이터 356건을 바탕으로, 2021년부터 2022년까지 화학물질안전원에서 제시한 화학사고 원인 조사 및 유사 사고 재발 방지 방안을 추가로 학습 데이터셋에 적용했다. 다층 퍼셉트론(Multi-Layer Perceptron) 모델의 경우 정확도 0.6590, 정밀도 0.6821로 분석되었고, 로지스틱 회귀(Logistic Regression) 모델은 정확도는 0.6647에서 0.7778로, 정밀도는 0.6790에서 0.7992로 향상되어 로지스틱 회귀 모델이 화학사고 원인 예측에 가장 효과적임을 확인하였다.

주요어 : 화학사고, 인공지능, 사고 원인 예측, 데이터 분석, 모델 성능 평가

Abstract This study aims to develop an artificial intelligence-based model for predicting the causes of chemical accidents, utilizing data on 865 chemical accident situations and cases provided by the Chemical Safety Agency under the Ministry of Environment from January 2014 to January 2024. The research involved training the data using six artificial intelligence models and compared evaluation metrics such as accuracy, precision, recall, and F1 score. Based on 356 chemical accident cases from 2020 to 2024, additional training data sets were applied using chemical accident cause investigations and similar accident prevention measures suggested by the Chemical Safety Agency from 2021 to 2022. Through this process, the Multi-Layer Perceptron (MLP) model showed an accuracy of 0.6590 and a precision of 0.6821. The Logistic Regression model improved its accuracy from 0.6647 to 0.7778 and its precision from 0.6790 to 0.7992, confirming that the Logistic Regression model is the most effective for predicting the causes of chemical accidents.

Key words : Chemical Accident, Artificial Intelligence, Accident Cause Prediction, Data Analysis, Model Performance Evaluation

*정희원, 한국공학대학교 IT반도체융합공학과 박사과정 (제1저자)

**정희원, 가톨릭관동대학교 책임연구원 (공동저자)

***준희원, 가톨릭관동대학교 의생명과학과 석사과정 (참여저자)

****정희원, 한국공학대학교 융합기술에너지대학원 교수 (교신저자)

*****정희원, 가톨릭관동대학교 의생명과학과 교수(교신저자)

접수일: 2024년 7월 7일, 수정완료일: 2024년 8월 20일

게재확정일: 2024년 9월 10일

Received: July 7, 2024 / Revised: August 20, 2024

Accepted: September 10, 2024

****Corresponding Author: kws@tukorea.ac.kr

Graduate school of Convergence Technology and Energy
Tech Univ of Korea

*****Corresponding Author: hjchoi@cku.ac.kr

Department of Biomedical Science, Catholic Kwandong
University

I. 서 론

화학사고는 화재나 폭발뿐만 아니라 누출의 위험성도 가지며, 독성과 부식성으로 인해 2차 피해 우려가 크다. 화학물질의 확산성, 비가시성, 유해성, 잔류성으로 인해 국민과 환경에 치명적인 영향을 미칠 수 있으며, 발생 유형이 복잡적이어서 대량 피해로 확대될 가능성이 높다.[1]

2024년 1월 3일 소방청 보도자료에 따르면, 2012년 9월 경북 구미의 한 공장에서 불산 가스 저장탱크가 폭발해 공장 근로자 5명이 사망하고 소방대원 등 18명이 부상을 입었다. 또한, 불산 가스가 인근 마을로 퍼져 3천여 명의 주민 등이 병원 진료를 받았으며, 주변 농작물과 가축도 피해를 입었다.

2018~2022년 전국 화학단지에서 발생한 화학물질 관련 사고는 145건으로, 사망자는 62명에 달했다. 국내 화학산업단지의 대부분이 노후시설로 인한 장비 결함, 저장탱크 부식, 관리 소홀 등으로 인해 화학 사고의 위험이 상존하고 있어 체계적이고 전문적인 대응이 요구된다. 소방청은 전국 주요 화학단지를 중심으로 전문 화학구조센터를 운영하고 있으나, 현장 관계자나 전문가를 통한 물질 확인에 오랜 시간이 걸리고, 그사이 피해가 확산할 우려가 있다. 이에 인공지능(Artificial Intelligence, AI)을 접목한 혁신적 대응 기술인 유해화학물질 관독에 인공지능 기술을 활용하여 초기 대응을 강화하고자 했다. 이에 따라, 2024년부터 실증사업으로 인공지능 활용, 재난 현장 영상 분석을 통해 유해화학물질 관독, 신속한 물질 정보 파악으로 초기 대응 강화 및 피해 범위 예측, 질산, 수산화나트륨, 시너, 등유 등 유해화학물질 총 10종의 학습데이터 생성 사업 등을 추진하고 있다.[2]

환경부 산하 화학물질안전원은 화학물질관리법 제48조(화학물질종합정보시스템 구축·운영)에 따라 화학물질 안전관리 정보, 화학사고 발생 이력 등과 관련된 정보를 제공하고 있으며 2021년도와 2022년도 화학 사고에 대한 원인 조사를 통해 유사 사고 재발 방지 방안을 제시하고, 이를 담은 사례집을 배포하였다.[3][4] 또한, 2025년까지 화학사고 발생 빈도가 높고 대형 사고 위험이 상존하는 전국 주요 노후 산단 15곳에 광화학 카메라, 인공지능 등 지능형 기술을 활용한 '원격 감시 체계'를 구축하여 운영할 예정이다.[5]

정부는 화학사고 관리에 인공지능 기술을 도입하고 있다. 인공지능 기술은 위험 예측, 유해화학물질 사고 대응, 건설 시 안전 대응 등 다양한 분야에서 활용되고 있으며, 이를 활용한 화학사고 예측, 관리, 통계, 분석 등 다양한 연구들이 진행되고 있다.

조철희(2017)는 119화학구조센터에서 대응한 출동 내역 및 사고 유형을 살펴보고, 화학 사고에 대한 대응 전략을 마련하고자 화학사고 통계 및 분석을 하였고, 2020년에는 국내 화학사고 조사 분석 및 효율적 대응 방안에 관해 연구하였다. [6][7]

김원덕(2021)는 기상 정보에 따른 화학사고 중대성 예측을 머신러닝 알고리즘 기술을 이용하여 연구를 진행하였고, 김남준(2022)는 인공지능 및 빅데이터를 활용한 화학사고 예측에 관한 연구를 진행하였다.[8][9]

이경현(2023)는 중대 산업사고 및 화학물질 사고 등을 예방하기 위한 인공지능 기반 통합 공정안전관리 시스템에 관한 연구를 수행했다.[10]

본 연구는 위의 선행 연구와 같이 화학물질 관련 사고 현황 및 사례 데이터를 활용한 인공지능 사고 원인 예측 모델에 관한 연구를 진행하고, 연구 결과를 바탕으로 사고 예방 및 대응 전략을 효율적으로 수립하여 안전한 작업 환경을 조성하는 데 기여하고자 한다.

II. 연구 개요

1. 학습데이터

인공지능 모델의 데이터 셋은 입력 특징(Feature)과 해당 입력에 대한 정답 레이블(Label)로 구성된다. 학습 데이터는 모델이 특정 입력에 대해 어떤 출력을 내야 하는지 배우도록 도와주고, 모델은 새로운 데이터에 대해 정확한 예측을 할 수 있도록 학습된다.

본 연구에서 사용된 학습데이터는 환경부 산하 화학물질안전원에서 제공하는 2014년 1월부터 2024년 1월까지 총 865건의 화학물질 사고 현황 및 사례 데이터이다. 표1에 865건의 사례 데이터를 사고 원인 및 사고 형태에 따라 분류하였다.[11]

표 1. 사고 원인 및 사고 형태의 집계 현황표[11]

Table 1. Summary of Accident Causes and Types

사고 원인	건수	사고 형태	건수
안전기준 미준수	363	누출	685
시설 결함	327	화재	66
운송 차량	167	폭발	69
자연재해	8	기타	45

사례 데이터에서 ‘사고 일자’, ‘사고 유형’, ‘사고 내용’, ‘제1 사고 물질’, ‘제2 사고 물질’, ‘제3 사고 물질’, ‘사고 원인’ 등의 항목을 인공지능 사고 원인 예측 모델의 학습데이터 셋으로 활용하였다.

연구에 활용된 6개의 인공지능 모델(랜덤 포레스트(Random Forest), 서포트 벡터 머신(Support Vector Machine, SVM), 로지스틱 회귀(Logistic Regression), 나이브 베이즈(Naive Bayes), K-최근접 이웃(K-Nearest Neighbors, KNN) 중에서 평가지표 결과의 값이 높은 인공지능 모델에는 화학물질안전원의 화학사고 원인 조사를 통해 유사 사고 재발 방지 방안으로 2021년도(5건), 2022년도(7건) 데이터를 학습데이터로 추가 적용하여 비교 분석하였다.[3][4]

2. 인공지능 모델

학습데이터를 사용하여 사고 원인을 예측하는 인공지능 모델을 개발하기 위해 총 6개의 인공지능 모델의 평가지표 값을 구하여 비교하였다. 본 연구에서 사용된 6개의 인공지능 분류 모델 관련 식은 다음과 같다.

2-1. 로지스틱 회귀 모델

로지스틱 회귀 모델은 이진 분류를 위한 선형 모델로서, 확률값을 출력하는 형태로 일반화된 선형 회귀이다. 다중 클래스 문제를 해결하기 위해 일반적으로 "일대다"(One vs Rest, OvR) 방식을 사용한다.

가. 시그모이드 함수(Sigmoid Function)

로지스틱 회귀 모델에서 출력값을 확률로 변환하기 위해 사용하는 함수로서 수식(1)과 같다.

$$\sigma(z) = \frac{1}{1 + e^{-x}} \quad \text{식(1)}$$

여기서 z 는 특징 벡터 \vec{x} 와 가중치 벡터 \vec{w} 의 선형 결합이다.

$$z = \vec{w}^T \vec{x} + b \quad \text{식(2)}$$

나. 손실 함수(Loss Function)

로지스틱 회귀 모델의 손실 함수로는 로그 손실 함수(Log loss) 또는 교차 엔트로피 손실 함수(Cross entropy loss)가 사용된다. 엔트로피 손실 함수는 모델의 예측 확률과 실제 레이블 간의 차이를 측정하여, 모델의 학습 과정에서 손실을 최소화하려고 한다.

$$L(y, \hat{y}) = -\frac{1}{m} \sum_{i=1}^m [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] \quad \text{식(3)}$$

여기서 y_i 는 실제 레이블(0 또는 1), \hat{y}_i 는 예측된 확률(모델이 예측한 해당 클래스일 확률), m 은 샘플의 수이다.

2-2. 다층 퍼셉트론 모델

신경망 모델의 일종으로, 입력 데이터를 여러 개의 층을 거쳐 비선형 변환하여 최종 출력값을 예측한다. 예측 과정에서 사용된 주요 수학적 개념과 식은 다음과 같다.

가. 입력과 가중치

식(4)는 다층 퍼셉트론의 뉴런에서 입력값과 가중치의 결합을 나타내는 식으로 z 는 선형 결합의 결과이다.

$$z = \vec{w}^T \vec{x} + b \quad \text{식(4)}$$

나. 활성화 함수

활성화 함수는 노드의 출력을 비선형 변환한다. MLP에서는 주로 ReLU(Rectified Linear Unit) 활성화 함수를 사용한다.

$$ReLU(z) = \max(0, z) \quad \text{식(5)}$$

다. 손실 함수 (Loss Function)

MLP 분류기는 교차 엔트로피 손실 함수를 사용하여 모델의 예측값을 실제 레이블과 비교한다.

$$L(y, \hat{y}) = -\sum_{i=1}^C y_i \log(\hat{y}_i) \quad \text{식(6)}$$

여기서 y_i 는 실제 레이블, \hat{y}_i 는 예측 확률, C 는 클래스의 수이다.

2-3. 랜덤 포레스트(Random Forest) 모델

다수의 결정 트리(Decision Tree)를 이용하여 예측을 수행하는 앙상블 학습 알고리즘이다. 각 트리는 무작위로 선택된 데이터와 특성으로 학습된다. 최종 예측은 모든 트리의 예측 결과를 투표 또는 평균 내는 방식으로 결정한다. 결정 트리는 데이터 분할 시 정보 이득이 최대가 되는 지점을 찾으며, 정보 이득은 분할 전후의 엔트로피(불확실성)의 감소량으로 정의된다.

Information Gain=Entropy(before)−Entropy(after)
 여기서 엔트로피는 다음과 같이 정의된다.

$$Entropy = - \sum_{i=1}^n p_i \log_2 p_i \quad \text{식(7)}$$

p_i 는 클래스 i 의 확률이다.

2-4. SVM 모델

SVM은 분류 문제를 해결하기 위한 지도 학습 모델로서 클래스 간의 최대 마진을 가진 초평면(하이퍼플레인)을 찾는 것과 데이터를 고차원 공간으로 매핑하여 비선형 분리를 가능하게 하는 두 가지 목적을 가진다.

가. 선형 SVM

선형 SVM은 두 클래스 사이의 최적의 초평면을 찾아 데이터를 분류한다. 최적의 초평면은 두 클래스 사이의 마진(두 클래스의 가장 가까운 데이터 포인트와 초평면 사이의 거리)을 최대화한다.

$$w \cdot x + b = 0 \quad \text{식(8)}$$

식(8)은 초평면의 방정식으로 w 는 초평면의 가중치 벡터, b 는 편향(term)을 나타낸다.

식(9)는 결정 경계로 두 클래스의 마진이 동일한 지점을 나타낸다.

$$\begin{cases} w \cdot x + b = 1 & \text{Positive Class} \\ w \cdot x + b = -1 & \text{Negative Class} \end{cases} \quad \text{식(9)}$$

나. 비선형 SVM

비선형 데이터를 분류하기 위해 SVM은 커널 트릭을 사용하여 데이터를 고차원 공간으로 매핑한다. 본 연구에서는 RBF 커널을 이용하였다.

RBF(Radial Basis Function) 커널

$$K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2) \quad \text{식(10)}$$

2-5. 다항 나이브 베이즈 모델

나이브 베이즈 분류기 중 하나로, 특히 텍스트 분류 문제에 자주 사용된다. 이 모델은 각 클래스의 단어 빈도 분포를 기반으로 텍스트를 분류한다.

가. 나이브 베이즈 정리

조건부 확률을 사용하여 클래스 C 가 주어진 데이터 포인트 x 에 속할 확률을 계산하고, 식(11)과 같은 형태를 가진다.

$$P(C|x) = \frac{P(x|C)P(C)}{P(x)} \quad \text{식(11)}$$

$P(C|x)$ 는 데이터 x 가 주어졌을 때 클래스 C 일

확률, $P(x|C)$ 는 클래스 C 가 주어졌을 때 데이터 x 일 확률, $P(C)$ 는 클래스 C 의 사전 확률, $P(x)$ 는 데이터 x 의 사전 확률이다.

나. 다항 나이브 베이즈

다항 나이브 베이즈 모델은 단어 빈도를 특징으로 하는 데이터에 적합하다. 각 클래스 C_k 에 대해 식(12)와 같은 확률을 계산한다.

여기서 x_i 는 텍스트의 i 번째 단어, $P(x_i|C_k)$ 는 클래스 C_k 에서 단어 x_i 의 확률이다.

$$P(C_k|x) = \frac{P(C_k) \prod_{i=1}^n P(x_i|C_k)}{P(x)} \quad \text{식(12)}$$

학습 단계로 모델은 주어진 데이터로부터 각 클래스 C_k 의 사전 확률 $P(C_k)$ 와 각 단어의 조건부 확률 $P(x_i|C_k)$ 를 추정한다.

$$P(C_k) = \frac{\text{Number of documents in class } C_k}{\text{Total number of documents}} \quad \text{식(13)}$$

$$P(x_i|C_k) = \frac{\text{Count of documents in class } C_k + \alpha}{\sum_{x_j \in V} (\text{Count of word } x_j \text{ in Class } C_k) + \alpha | V|} \quad \text{식(14)}$$

여기서 α 는 라플라스 스무딩(Laplace Smoothing) 파라미터, V 는 단어의 전체 집합이다.

2-6. K-최근접 이웃 모델

비모수적 방법의 지도 학습 알고리즘으로, 새로운 데이터 포인트를 기존 데이터와 가장 가까운 K개의 이웃으로 분류한다.

가. 거리 측정

KNN 알고리즘은 주로 유클리드 거리를 사용하여 포인트 간의 거리를 측정한다. 유클리드 거리(Euclidean Distance)는 식(15)와 같이 정의된다.

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad \text{식(15)}$$

x 와 y 는 두 데이터 포인트, n 은 특징의 수이다.

나. 다수결 투표

새로운 데이터 포인트의 클래스를 결정하기 위해 K개의 최근접 이웃의 클래스를 투표하여 다수결로 결정한다.

$$\hat{y} = \operatorname{argmax}_c \sum_{i \in N_k} \Pi(y_i = c) \quad \text{식(16)}$$

식(16)에서 N_k 는 K개의 최근접 이웃, Π 는 인디케이터 함수로, 이웃의 클래스 y_i 가 특정 클래스 c 와 일치하면 1, 그렇지 않으면 0을 반환한다.

3. 데이터 전처리

데이터 전처리는 모델이 학습할 수 있도록 원시 데이터를 준비하는 과정으로 데이터의 품질을 향상시키고 모델의 성능을 극대화하기 위해 필수적인 단계이다. 데이터 전처리에는 여러 단계가 포함될 수 있으며, 여기에는 결측값 처리, 데이터 정규화, 범주형 데이터 인코딩, 특성 선택 및 추출 등이 포함된다.

가. 텍스트 데이터 벡터화

TF-IDF(Term Frequency-Inverse Document Frequency) 벡터화 방법을 사용하여 텍스트 데이터를 수치화한다.

$$TF-IDF(t, d) = TF(t, d) \times IDF(t) \quad \text{식(17)}$$

식(17)에서 $TF(t, d)$ 는 문서 d 에서 단어 t 의 빈도, $IDF(t)$ 는 단어 t 의 역문서 빈도이다.

나. 범주형 데이터 라벨 인코딩

라벨 인코딩을 통해 범주형 데이터를 수치형 데이터로 변환한다.

다. 표준화

표준화는 데이터의 분포를 평균이 0, 표준편차가 1이 되도록 변환한다.

$$X_{scaled} = \frac{X - \mu}{\sigma} \quad \text{식(18)}$$

식(18)에서 μ 는 평균, σ 는 표준편차이다.

4. 평가지표

인공지능 모델의 성능을 평가하기 위해 다양한 평가 지표가 사용되며, 이 평가지표들은 모델의 예측 성능을 여러 측면에서 평가하며, 모델의 강점과 약점을 파악하는 데 도움을 준다. 본 연구에서 사용되는 평가지표들을 아래에 나타내었다.

4-1. 정확도(Accuracy)

전체 예측 중에서 올바르게 분류된 샘플의 비율을 말한다.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad \text{식(19)}$$

여기서 TP(True Positive)는 실제 Positive인 샘플을 Positive로 올바르게 분류한 경우의 수

TN(True Negative)는 실제 Negative인 샘플을 Negative로 올바르게 분류한 경우의 수

FP(False Positive)는 실제 Negative인 샘플을 Positive로 잘못 분류한 경우의 수

FN(False Negative)는 실제 Positive인 샘플을 Negative로 잘못 분류한 경우의 수를 나타낸다.

4-2. 정밀도(Precision)

Positive라고 예측한 샘플 중 실제로 Positive인 샘플의 비율을 나타낸다.

$$Precision = \frac{TP}{TP + FP} \quad \text{식(20)}$$

4-3. 재현율(Recall)

실제로 Positive인 샘플 중 모델이 Positive로 예측한 샘플의 비율을 나타낸다.

$$Recall = \frac{TP}{TP + FN} \quad \text{식(21)}$$

4-4. F1 점수 (F1 Score)

정밀도(Precision)와 재현율(Recall)의 조화 평균으로, 두 메트릭의 균형을 평가하는 데 사용된다. 불균형한 클래스 분포에서 모델의 성능을 종합적으로 평가하는 데 유용하다. F1 점수는 모델이 예측한 Positive 클래스 중 실제로 Positive인 비율(정밀도)과 실제로 Positive인 샘플 중 모델이 Positive로 예측한 비율(재현율)의 조화를 반영한다. 수식은 식(22)와 같다.

$$F1\ Score = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad \text{식(22)}$$

F1 점수는 0에서 1 사이의 값을 가지며, 1에 가까울수록 모델의 성능이 좋음을 나타낸다. 특히, 클래스 불균형이 심한 경우, F1 점수는 모델 성능을 더 공정하게 평가할 수 있는 지표로 사용된다.

4-5. 혼동 행렬(Confusion Matrix)

모델의 예측 결과를 실제 클래스와 비교하여 TP, FP, FN, TN을 표현하는 행렬이다. 모델의 예측 오류를 시각적으로 확인할 수 있다.

혼동 행렬은 분류 모델의 성능을 평가하는 데 사용되는 표이며, 실제 클래스와 모델이 예측한 클래스를 비교하여 구성된다. 주로 이진 분류(Binary Classification)에서 사용되지만, 다중 클래스(Multi Class) 분류에서도 적용할 수 있다.

4-6. ROC(Receiver Operating Characteristic) Curve 및 AUC (Area Under the Curve)

ROC Curve는 다양한 임계값에서 분류 모델의 성능을 평가하는 그래프로 y축은 TPR(True Positive Rate), x축은 FPR(False Positive Rate)을 나타낸다.

$$TPR = \frac{True\ Positives}{True\ Positives + False\ Negatives} \quad \text{식(23)}$$

$$FPR = \frac{False\ Positives}{False\ Positives + True\ Negatives} \quad \text{식(24)}$$

AUC는 ROC 곡선 아래의 면적을 나타내며, 0.5에서 1 사이의 값을 가진다. AUC가 1에 가까울수록 모델의 성능이 좋을 나타낸다.

사고 원인 예측 모델에 사용된 평가지표와 ROC 곡선은 모델의 성능을 종합적으로 평가하는 데 중요한 역할을 하며, 모델의 각 클래스에 대해 얼마나 잘 예측하는지, 오류가 얼마나 발생하는지를 시각적으로 확인할 수 있다.

III. 본 문

본 연구는 6개 사고 원인 예측 인공지능 분류 모델을 사용했으며, 사고 원인 예측 인공지능 분류 모델 프로그램의 순서도를 그림1에 정리하였다.

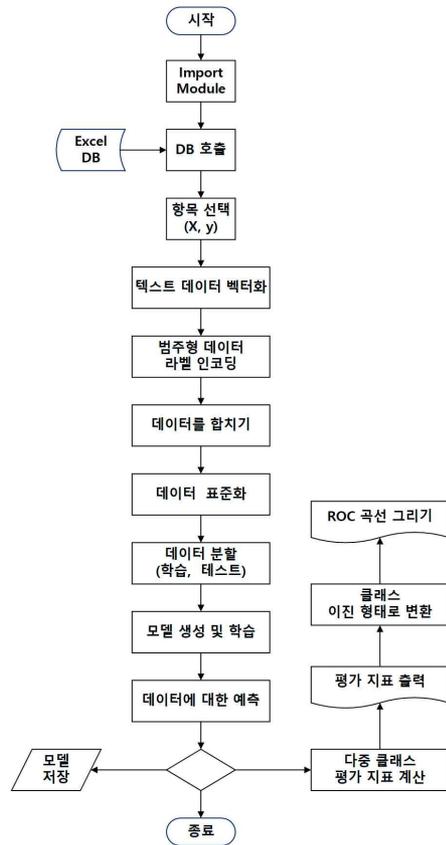


그림 1. 사고 원인 예측 인공지능 분류 모델 프로그램 순서도
Figure 1. Artificial intelligence predicts the cause of the accident and the scheduled time

위의 순서도를 설명하면, 다음과 같다.

가. 데이터 수집 및 준비

Excel로 처리된 데이터베이스를 불러온 후, '사고 원인'을 y 변수로 설정하고, x 항목으로는 '사고 일자', '사고 유형', '사고 내용', '제1 사고 물질', '제2 사고 물질', '제3 사고 물질' 등을 선택하였다.

표 2. 인공지능 모델별 평가지표 결과 값
Table 2. Evaluation index results for each artificial intelligence model

평가지표	나이브 베이즈	K-최근접 이웃	서포트 벡터 머신	랜덤 포레스트	로지스틱 회귀	다층 퍼셉트론
Accuracy	0.3526	0.4335	0.5954	0.6474	0.6647	0.6590
Precision	0.4267	0.5519	0.6584	0.6614	0.6790	0.6821
Recall	0.3526	0.4335	0.5954	0.6474	0.6647	0.6590
F1 Score	0.3769	0.3711	0.5885	0.6490	0.6662	0.6632
Confusion Matrix	$\begin{pmatrix} 21 & 27 & 7 & 14 \\ 17 & 32 & 10 & 10 \\ 2 & 21 & 8 & 3 \\ 1 & 0 & 0 & 0 \end{pmatrix}$	$\begin{pmatrix} 20 & 49 & 0 & 0 \\ 15 & 54 & 0 & 0 \\ 8 & 25 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}$	$\begin{pmatrix} 38 & 31 & 0 & 0 \\ 17 & 52 & 0 & 0 \\ 6 & 15 & 13 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}$	$\begin{pmatrix} 44 & 24 & 1 & 0 \\ 23 & 44 & 2 & 0 \\ 4 & 6 & 24 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}$	$\begin{pmatrix} 42 & 25 & 2 & 0 \\ 21 & 48 & 0 & 0 \\ 6 & 3 & 25 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}$	$\begin{pmatrix} 36 & 28 & 2 & 3 \\ 17 & 50 & 1 & 1 \\ 2 & 4 & 27 & 1 \\ 0 & 0 & 0 & 1 \end{pmatrix}$

나. 데이터 전처리

텍스트 데이터를 벡터화하고, 범주형 데이터는 라벨 인코딩을 통해 수치형 데이터로 변환하였다. 또한, 데이터 합치기 및 표준화를 통해 데이터의 분포를 평균이 0 이고 표준편차가 1이 되도록 변환하였다.

다. 모델 학습 및 평가

데이터를 학습용과 테스트용으로 8 : 2 비율로 분할 하고, 6개의 인공지능 모델을 생성하고 학습하도록 하였다. 학습된 모델을 사용하여 데이터를 예측하고, 평가 지표를 계산하여 모델의 성능을 평가하였다.

라. 모델 저장

모델 성능을 평가하여 기준값 이상의 결과 값을 갖는 사고 원인 예측 인공지능 모델을 저장하도록 하였다.

마. 모델 성능 비교

6개의 인공지능 모델(K-최근접 이웃, 서포트 벡터 머신, 나이브 베이즈, 랜덤 포레스트, 로지스틱 회귀, 다 층 퍼셉트론)에 대해 다중 클래스 평가지표인 정확도, 정밀도, 재현율, F1 점수를 계산하여 성능을 평가하였다. 또한, ROC 곡선을 그려 각 모델의 성능을 시각적으로 평가할 수 있도록 하였다.

각 인공지능 모델의 실행 결과를 표2)에 정리하였다. 평가지표 결과 값이 큰 다층 퍼셉트론 모델, 로지스틱 회귀 모델을 이용하여 2020년 1월부터 2024년 1월까지의 화학물질 사고 현황 및 사례 데이터 356건과 화학물질 안전원에서 제시한 화학사고 원인 조사와 유사 사고 재발 방지 방안 제시안(12건, 2021~2022년)을 학습데이터에 추가 적용한 결과 값을 표3)와 표4)에 정리하였다.

표3)과 표4)의 I 은 2014년부터 2024년까지의 화학물질 사고 데이터 865건에 대한 평가지표 결과 값을, II 은 2020년 1월부터 2024년 1월까지의 사고 데이터 356 건에 대한 평가지표 결과 값을, III 은 II의 사고 데이터 356건에 2021~2022년 화학물질안전원이 제시한 화학 사고 원인 조사와 유사사고 재발 방지 방안 12건의 데이터를 추가 적용한 학습데이터를 사용하여 나온 평가 지표 결과 값을 정리하였다.

실행 결과, 표3) 의 다층 퍼셉트론 모델에서는 2021~2022년 화학사고 원인 조사와 재발 방지 방안을 추가한 데이터 셋에서 정확도와 정밀도, 재현율, F1 점수에 있어

성능이 저하되었다. 추가된 데이터가 모델의 일반화 성능에 부정적인 영향을 미쳤을 가능성이 있다고 판단하였다.

표 3. 다층 퍼셉트론 모델 평가지표 비교표

Table 3. Multilayer perceptron neural network evaluation index comparison table

평가지표	I	II	III
Accuracy	0.6590	0.6944	0.625
Precision	0.6821	0.7293	0.6467
Recall	0.6590	0.6944	0.625
F1 Score	0.6632	0.6945	0.6287
Confusion Matrix	$\begin{pmatrix} 36 & 28 & 2 & 3 \\ 17 & 50 & 1 & 1 \\ 2 & 4 & 27 & 1 \\ 0 & 0 & 0 & 1 \end{pmatrix}$	$\begin{pmatrix} 17 & 4 & 0 \\ 12 & 23 & 4 \\ 1 & 1 & 10 \end{pmatrix}$	$\begin{pmatrix} 12 & 7 & 2 \\ 9 & 24 & 6 \\ 2 & 1 & 9 \end{pmatrix}$

표 4. 로지스틱 회귀 모델 평가지표 비교표

Table 4. Logistic regression evaluation index comparison table

평가지표	I	II	III
Accuracy	0.6647	0.75	0.7778
Precision	0.6790	0.7655	0.7992
Recall	0.6647	0.75	0.7778
F1 Score	0.6662	0.7463	0.7707
Confusion Matrix	$\begin{pmatrix} 42 & 25 & 2 & 0 \\ 21 & 48 & 0 & 0 \\ 6 & 3 & 25 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}$	$\begin{pmatrix} 13 & 8 & 0 \\ 5 & 34 & 0 \\ 2 & 3 & 7 \end{pmatrix}$	$\begin{pmatrix} 13 & 8 & 0 \\ 3 & 36 & 0 \\ 2 & 3 & 7 \end{pmatrix}$

표4) 의 로지스틱 회귀 모델 실행 결과에서 II의 결과 값의 정확도와 정밀도, 재현율, F1 점수 모두에서 I의 결과에 비하여 13% 이상의 성능향상을 확인하였고, 화학물질안전원이 제시한 화학사고 원인 조사와 유사 사고 재발 방지 방안이 추가 적용된 III의 결과 값은 I의 결과 값에 비하여 정확도와 정밀도, 재현율, F1 점수 모두에서 약 17%의 성능 개선을 확인할수 있었다.

추가된 데이터가 모델의 성능 향상에 기여하고 있음을 판단하였다.

그림2, 그림3, 그림4은 로지스틱 회귀 모델의 I, II, III의 학습데이터 대한 ROC커브 및 AUC 값에 대한 것이다.

로지스틱 회귀 모델의 ROC 커브 또한 각 클래스에 대한 TPR과 FPR를 나타내며, AUC가 1에 가까울수록 분류 성능이 우수하다. 그림2에서 class는 사고 원인을 나타내며 class 0(■)은 시설 결함, class 1(▲)은 안전 기준 미준수, class 2(★)는 운송 차량, class 3(●)은 자연재해 등의 결과 값이다. 그림3과 그림4에서 class 0 (■)은 시설 결함, class 1(▲)은 안전기준 미준수, class 2(★)는 운송 차량 등의 결과 값을 나타낸다.

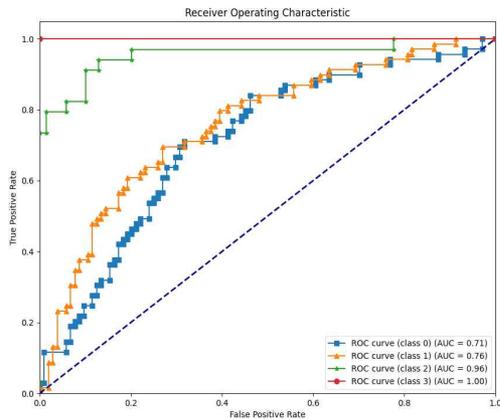


그림 2. 로지스틱 회귀(865건) ROC Curve
Figure 2. Logistic regression (865 cases) ROC Curve

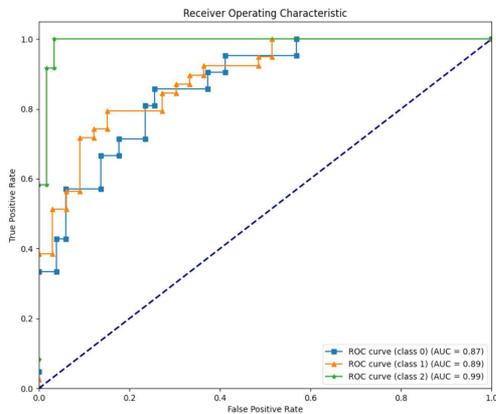


그림 3. 로지스틱 회귀(356건) ROC Curve
Figure 3. Logistic regression (356 cases) ROC Curve

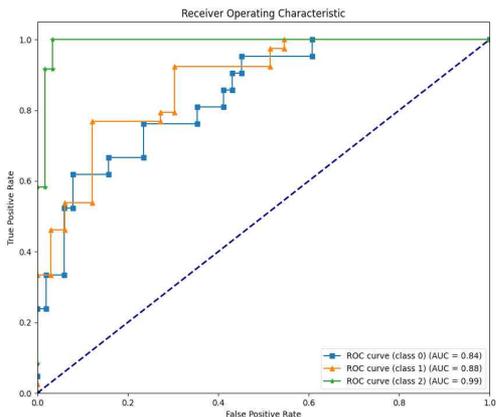


그림 4. 로지스틱 회귀(356건+화학사고 내용, 원인 재작성) ROC Curve
Figure 4. Logistic regression (356 cases + chemical accident details and causes rewritten) ROC Curve

ROC 커브와 AUC를 통해 각 클래스에서의 예측 성능을 시각적으로 평가할 수 있으며, III 데이터 셋에서 로지스틱 회귀 모델 AUC 값이 class 0은 0.84, class 1은 0.88, class 3은 0.99로, 분류 성능이 우수함을 알 수 있다.

V. 결 론

본 연구에서는 2014년부터 2024년까지의 화학사고 데이터를 바탕으로 6개의 인공지능 모델을 사용하여 사고 원인 예측 모델을 개발하였다. 그 결과, 로지스틱 회귀와 다층 퍼셉트론 모델이 가장 높은 성능을 보였으며, 특히 2020년부터 2024년까지의 데이터를 추가 학습 데이터로 사용하여 모델 성능이 향상됨을 확인하였다.

다층 퍼셉트론 모델은 정확도 0.6590, 정밀도 0.6821로 분석되었고, 로지스틱 회귀 모델의 경우 정확도가 0.6647에서 0.7778로 향상되었으며, 정밀도는 0.6790에서 0.7992로 증가하였다. 추가된 데이터로 인해 모델의 성능이 향상되었음을 보여준다.

결론적으로, 로지스틱 회귀 모델이 화학사고 원인 예측에 가장 적합한 모델임을 확인할 수 있었고, 다층 퍼셉트론 모델은 일부 클래스에서 예측 정확도가 낮았지만, 전체적인 성능은 여전히 우수하였다.

향후 연구로는 다양한 인공지능 모델을 결합하거나, 같은 모델의 다른 하이퍼파라미터 조합을 사용하여 앙상블 기법의 적용이 가능하며, 화학사고 원인을 더 정확하게 예측하고, 실시간 경보 시스템을 구축하여 안전 사고를 예방할 수 있을 것으로 예상된다. 이러한 연구를 통해 안전한 작업 환경을 조성하고, 사고 예방 및 대응 전략을 효율적 수립이 가능할 것이다.

References

- [1] Ji-sun You, & Yeong-Jin Chung, "Case Analysis of the Harmful Chemical Substances Spill", Journal of Korean Institute of Fire Science & Engineering, 2014, vol.28 no.6, 90-98.
- [2] Hyun-joo Moon "Expansion of the scope of first aid work for 119 paramedics...Expected to improve survival rate through rapid treatment of seriously ill patients", National Fire Agency, 2024.01.03. https://www.nfa.go.kr/nfa/news/pressrselease/press/?boardId=bbs_000000000000010&mod

- e=view&cntId=2060&category=&pageIdx=13&searchCondition=all&searchKeyword
- [3] Chun-hwa Park “Propose measures to prevent recurrence of similar accidents by investigating the causes of chemical accidents”, Ministry of Environment, 2021.12.24. <https://www.me.go.kr/home/web/board/read.do?pagerOffset=0&maxPageItems=10&maxIndexPages=10&searchKey=title&searchValue=&menuId=10525&orgCd=&boardId=1497470&boardMasterId=1&boardCategoryId=39&decorator>
- [4] Byung-Hoon Kim, “Suggestion of measures to prevent accidents by investigating the cause of chemical accidents”, Ministry of Environment, 2022.12.29. <https://www.me.go.kr/home/web/board/read.do?pagerOffset=0&maxPageItems=10&maxIndexPages=10&searchKey=title&searchValue=&menuId=10525&orgCd=&condition.fromDate=2022-12-01&condition.toDate=2022-12-31&boardId=1571180&boardMasterId=1&boardCategoryId>
- [5] Yong-soon Lim “Ministry of Environment establishes safety net to monitor chemical leaks in Incheon Namdong National Industrial Complex”, Ministry of Environment, 2022.03.21. <https://www.me.go.kr/home/web/board/read.do?pagerOffset=0&maxPageItems=10&maxIndexPages=10&menuId=10525&orgCd=&boardId=1514940&boardMasterId=1&boardCategoryId=&decorator>
- [6] Chul-Hee Cho, Dong Won Lee, & Sung Yeon Kim “A Review of Statistics & Analysis on the Chemical Accidents in Korea 2017 : Focus on the NFA, National 119 Rescue Headquarters” Korean Journal of Hazardous Materials, 2018, vol.6, no.1, 37-46,
- [7] Chul-Hee Cho, Sin-Woong Choi, Sang-Hee Lee, Jung In Kim & Tae-Woo Kim. “Study on the Chemical Accidents Investigation and Effective Response System in Korea 2020” Korean Journal of Hazardous Materials 2021, vol.9, no.2, 68–75.
- [8] Won-deok Kim, “A prediction of the severity of chemical accidents based on weather information using machine learning algorithm technology.” master degree, SEOUL NATIONAL UNIVERSITY OF SCIENCE AND TECHNOLOGY, 2021
- [9] Nam-Jun Kim, “Prediction of chemical accidents using artificial intelligence(AI) and big data (Analysis of the effects of meteorological conditions and material characteristics on the occurrence of chemical accidents).” doctoral degree, TECH UNIVERSITY OF KOREA, 2022
- [10] Kyung-Hyun Lee, Rack-June Baek, & WooSu Kim. “A Study on Applying Novel Reverse N-Gram for Construction of Natural Language Processing Dictionary for Healthcare Big Data Analysis.” The Journal of Convergence on Culture Technology. 2024, vol.10, no.3, 391–396
- [11] National Institute of Chemical Safety. <https://icis.me.go.kr/pageLink.do>
- [12] Hyo-Jung Oh, & Bo-Hyun Yun. “Valid Data Conditions and Discrimination for Machine Learning: Case Study on Dataset in the Public Data Portal.” Journal of Internet of Things and Convergence. 2022, vol.8, no.1, 37 - 43

※ This work was supported by Regional Innovation Strategy (RIS) through the National Research Foundation of Korea (NRF) funded by the Ministry of Education (MOE) (2022RIS-005).