

이슈보고서 26-01

ISSUE
REPORT

자본시장 심리지수 시리즈 1

자본시장 심리지수의 구축과 활용

노성호

자본시장 심리지수의 구축과 활용

저자 노성호*

본 연구는 '자본시장 심리지수(Capital Market Sentiment Index: CMSI)'의 개념적인 배경과 구축 방법, 통계적 특성 및 활용 방안을 종합적으로 논의한다. 거시경제 불확실성과 시장 변동성이 확대되는 상황에서 정량적 지표에 더하여 비정형 정보를 활용할 필요성이 증가하고 있다. 이와 같은 배경에서 자본시장 심리지수는 자본시장 참여자들의 정성적 판단, 감정, 기대와 우려를 정량적으로 측정하도록 고안되었다. 특히, 트랜스포머(transformer) 기반 범용 대형언어모형(Large Language Model: LLM)을 활용하여 국내 뉴스 데이터를 학습하고 이를 바탕으로 시장 참여자들의 관점을 정량화한 지표를 구축하는 방법을 본 연구에서 상세하게 제시한다.

자본시장 심리지수는 복잡한 LLM에 기반하고 있으면서도 지수의 투명성과 설명가능성을 제고하였다는 점에서 의미가 있다. 우선, 뉴스 텍스트의 전처리 과정 및 학습 알고리즘을 상세하게 설명하여 모형의 투명성과 재현 가능성을 강화하였다. 더불어 샵플리 값(Shapley value) 분석을 활용하여 모형의 분석 결과를 입력된 단어별 기여도로 분해하여 정량적인 해석을 가능하게 하였다. 문장별로 해석된 시장 심리지표는 일별 평균과 분산 값을 기초로 다층적인 심리지수 시계열 자료를 구축하고 이들의 통계적 특성을 바탕으로 활용 방안을 제시하였다.

자본시장 심리지수는 근거기반(evidence-based) 의사결정을 보조할 수 있다는 점에서 유용성을 찾을 수 있다. 특히, 비정형 문자 정보에 기반한 고빈도 지표로서 투자자 심리의 실시간 모니터링 및 정책 효과 검증 등의 영역에서 활용 가치가 높다. 나아가, 향후 다차원 분류체계 도입, 멀티모달 모형 활용 등을 통하여 지수 체계를 고도화하는 연구의 기반이 될 것으로 기대한다.

* 본고의 견해와 주장은 필자 개인의 것이며, 자본시장연구원의 공식적인 견해가 아님을 밝힙니다.

금융산업실 연구위원 노성호 (nohs@kcmi.re.kr)

** 발행: 2026년 1월 5일

1. 서론

1. 연구의 배경 및 목적

글로벌 경제의 불확실성은 지속되는 인플레이션 압력, 금리 및 환율 변동, 지정학적 리스크와 공급망 위험 등으로 인하여 크게 높아지고 있다. 이와 같이 불확실성이 높은 시기에는 심리적 요인이 자산 가격에 영향을 미칠 가능성이 높아 이를 적시에 포착할 수 있는 지표가 유용하다(Huang et al., 2019). 특히, 투자자 심리는 시장의 과민 반응 또는 진정 국면을 미리 알려주는 예측 변수로서 기능할 수 있다(Baker & Wurgler, 2007; Huang et al., 2015)는 점에서 동행 또는 선행 지표로서 정보적 유용성이 크다. 나아가 시장의 정보 비대칭성과 그로 인한 심리적 요인은 경제의 펀더멘털(fundamental)에 대한 신호로 작용하여 경기 변동에도 영향을 미칠 수 있어(Benhabib et al., 2016) 이에 대한 정량화된 지표를 구축하는 것은 투자 및 정책 의사결정의 효율성과 적합성을 제고할 수 있어 중요한 과제이다.

이에 본 연구에서는 주식시장에 대한 뉴스를 학습한 대형언어모형(Large Language Model: LLM)을 바탕으로 자본시장 심리지수(Capital Market Sentiment Index: CMSI)를 구축하고 이를 활용하는 방안을 제시하였다. 한국어 문서 분류에 강점을 지닌 범용 LLM을 바탕으로 증권 시장에 대한 뉴스를 학습하여 투자 심리의 변화를 이해하고 정량화하는데 특화된 모형을 제안한다. 이를 위하여 대규모의 텍스트 데이터를 체계적으로 수집하고 처리하는 과정을 상세히 설명하였다. 나아가 모형을 기반으로 구축한 지수의 의미를 해석하고 정보적인 유의성을 주가 수익률과 변동성에 영향을 미치는 사건을 중심으로 분석하는 것을 목적으로 한다.

자본시장 심리지수의 배경 이론과 활용 방안은 세 편의 보고서로 나누어서 상세하게 논의한다. 본 보고서는 세 편의 연작 중 첫 번째로 CMSI의 구축 과정을 구체적으로 소개하고 산출된 지수의 통계적인 특징을 분석한다. 두 번째 보고서로 장보성(2026)은 장기 시계열 자료를 바탕으로 CMSI와 주요 금융 및 거시경제 지표와의 인과관계를 분석하였다. 이를 통해 CMSI가 기존의 정량적 지표와 차별화되어 주식시장의 동향에 대한 유의성 높은 정보를 추가적으로 제공하고 있음을 보인다. 세 번째로 김민기(2026)는 CMSI를 활용하여 추정된 시장 참여자의 (투자) 심리가 종목별 수익률에 이질적인 영향이 있음을 확인하였다. 이는 정성적인 요인이 주요 시장 활동(예: IPO)에 미치는 영향을 실증적으로 보여준 사례로서 CMSI가 투자 및 정책 의사결정에 유용한 지표임을 정량적으로 확인할 수 있다.

2. 관련 선행 연구 및 본 연구의 특징

머신러닝(machine learning) 및 인공지능(Artificial Intelligence: AI) 방법론의 발전에 힘입어 빅 데이터를 활용하여 경제 전반 또는 금융시장 동향에 대한 시의성 있는 정보를 추출하려는 시도가 다방면으로 진행되고 있다. 특히, 체계적으로 수집된 뉴스는 시장 전반에 대한 대중의 인식을 반영하고 있다는 점에서 정보의 시의성과 포괄성이 높은 자료이며 양적으로도 LLM을 학습시키기 적합한 자료이다.

뉴스 데이터를 학습시킨 모형은 개발 목적에 따라 크게 거시 경기 변동을 추적하는 모형과 금융 정보의 추출에 특화된 모형으로 나누어 볼 수 있다. 우선, 전자에 해당하는 경우에 국내 뉴스를 활용한 한글 기반 모형으로는 단어별 비중을 분석한 김현중 외(2019), 트랜스포머(transformer) 기반 언어모형⁰¹을 활용하여 전반적인 경기 변동을 나타내는 지표를 개발한 서범석 · 이영환 · 조영배(2022)가 그 예시이다. 해외 뉴스를 활용한 예시로 Huang et al.(2019)은 20개 국가의 뉴스를 포괄적으로 분석하여 글로벌 경기 변동 양상을 추적하는 지수를 개발하였다. Bybee et al.(2024)은 대량의 Wall Street Journal 기사를 학습시킨 모형을 바탕으로 구축한 뉴스 기반 지표와 기존의 경제 및 시장 활동 지표와의 장기적인 연관관계를 분석하였다.

금융 정보에 특화된 연구로는 텍스트 데이터에서 자산 가격 결정 요인을 추출하려는 시도를 들 수 있다. Da et al.(2015)은 인터넷 검색 기록을 활용하여 시장 전반에 대한 투자 심리를 정량화 하였으며, Bybee et al.(2023)은 뉴스 기사를 주제별로 분류하여 자산 가격 결정 요인을 추출하였다. Kim et al.(2024)은 사전 학습된 LLM을 국내 증권 뉴스를 활용하여 시장 심리를 분석하는 성능을 강화하였다는 점에서 본 연구와 가장 밀접하게 연관되어 있다. 해당 연구에서는 구축된 LLM을 바탕으로 정량화된 지수를 계산하고 이를 변동성, 위험회피도 등과 비교하여 지수의 유용성을 입증한 것으로 평가된다.

본 연구는 이와 같은 방법론을 국내 뉴스 데이터에 적용하여 한국 증시에 대한 시장 참여자의 인식을 정량화하고 동향을 적시에 파악할 수 있는 이론적인 기반을 구축한다는 점에서 의의를 지닌다. 자본시장 심리지수는 문자로 표현된 비정형 데이터에서 시장 동향에 유의미한 정보를 추출한다는 점에서 기존의 시장 지표와 차별화된 정보적 가치를 제공하고 있다. 나아가 복잡한

01 트랜스포머(transformer)는 중첩된 신경망(neural network) 모형의 일종으로 글을 인식할 때 문맥 안에서 각 단어가 내포하는 의미를 함께 고려할 수 있도록 하여 언어모형의 대형화에 기여한 기술로 평가된다. 트랜스포머 모형의 도입으로 촉발된 언어모형의 대형화 기초에 대한 설명은 노성호(2024)의 논의를 참고하기 바란다.

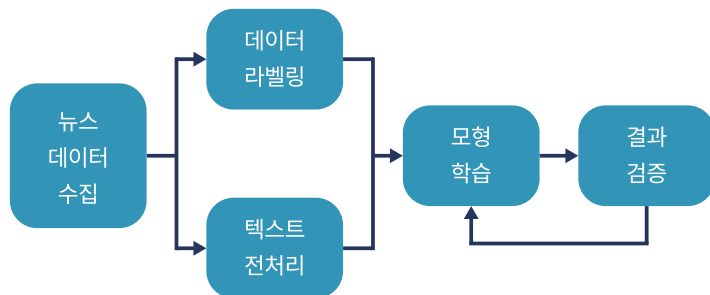
LLM에 기반하여 측정된 지표의 계산 절차를 구체적으로 설명하고 학습 데이터를 구성하는 텍스트를 문장 및 형태소 단위로 분해하고 기여도를 계산하여 다소 추상적일 수 있는 지수값에 대하여 인간의 관점에서 해석이 가능한 단어에 기반하여 해석하는 방안을 제시하였다. 이와 같은 논의는 흔히 블랙박스(black box)로 일컬어지는 초거대 모형에 기반하여 계산된 정량적 지표의 투명성과 설명가능성을 제고하였다는 점에서 의의가 있다.

II. 자본시장 심리지수 기반 모형

1. 뉴스 데이터의 학습 과정

CMSI의 계산을 위하여 본 연구에서는 한국어 LLM에 기반하여 미세조정(fine-tuning)을 거쳐 뉴스 기사에 포함된 문장의 특성을 분류할 수 있는 모형을 구축하였다. 우선, 기반이 되는 모형은 한국어 자료를 대량으로 학습하여 한글 문서 분류에 우수한 성능을 보인 KLUE-BERT를 사용하였으며⁰², 미세조정 절차는 Kim et al.(2024)을 참고하였다.

<그림 II-1> 자본시장 심리지수 기반모형 구축 과정



자료: 저자 작성

기반 모형을 미세조정하는 과정은 크게 <그림 II-1>과 같은 흐름으로 구조화할 수 있다. 우선, 학습의 대상인 뉴스 데이터를 수집하고 이를 분류하는 과정이 선행되어야 한다. 더불어 텍스트 전처리(pre-processing) 과정을 거쳐 한글 문장을 머신러닝 알고리즘에 입력가능한 형태로

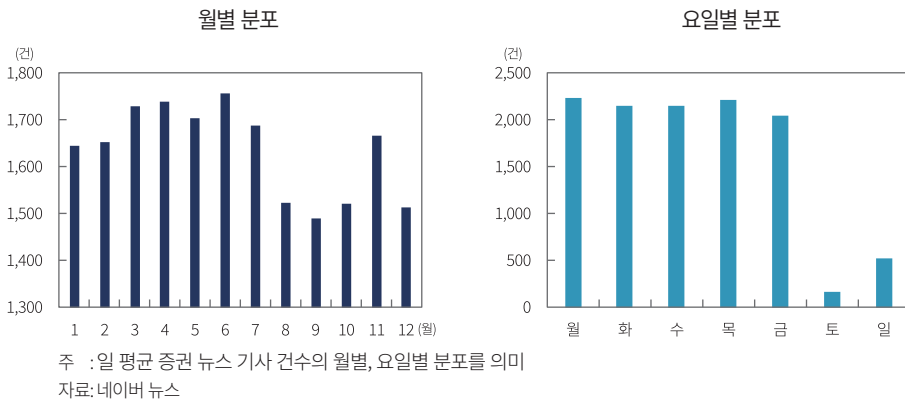
02 KLUE-BERT의 특징과 성능은 Park et al.(2021)을 참조하기 바란다.

가공하는 절차를 거친다. 이와 같이 준비된 뉴스 데이터는 정해진 알고리즘에 따라 기반 모형을 (재)학습시키는데 사용되며 마지막으로 학습된 모형의 텍스트 분류 성능을 평가하여 필요시 모형을 재조정하는 과정을 거치게 된다. 각 단계에서 수행되는 작업에 대하여 보다 자세한 설명은 아래와 같다.

가. 학습 데이터 개요

모형의 학습 및 지수의 생산에 이르는 과정은 크게 데이터 수집과 전처리, 데이터 라벨링, 분류 모형의 학습과 검증 단계로 나눌 수 있다. 이 중 첫 번째 단계는 모형의 학습에 사용될 뉴스 기사의 수집이다. 대규모의 기사 원문 자료는 인터넷에 공개된 뉴스 기사 텍스트에 대한 웹크롤링(web crawling) 방식으로 수집하였다. 특히 주식시장에 대한 정보만으로 한정하기 위하여 여러 분야 중 증권 시장에 대한 뉴스로 한정하여 학습 데이터를 구축하였다.

<그림 II-2> 학습 데이터 분포



수집된 데이터는 2010년 1월 1일부터 2025년 3월 31일에 걸쳐 일 평균 1,600여 건의 증권 뉴스를 포함한다. 수집된 데이터의 월별, 요일별 분포는 <그림 II-2>에 나타나 있다. 월별로는 3월에서 7월 사이에 상대적으로 많은 데이터가 분포하고 있다. 요일별로는 주말에 공개된 기사가 유의하게 적은 것을 고려하여 시장 지표와의 상관관계 등 향후 분석에서 휴일을 제외한 거래일로 표본을 한정하여 분석을 수행하였다.

나. 데이터 라벨링

지수를 계산하는 기반 모형은 인간 연구진에 의해 검증된 데이터를 사용하여 지도학습(supervised learning) 방식으로 구축되며 이를 위해서는 원문의 일부를 연구진이 직접 분류하는 라벨링 (labelling) 작업이 선행되어야 한다. 이와 같은 방법론을 선택한 배경에는 원문을 분류하지 않고 그대로 사용하는 비지도학습(unsupervised learning)과 비교하여 지도학습이 가진 두 가지 장점을 고려하였다. 우선, 뉴스 텍스트에 내재된 복합적인 문맥, 어조, 은유 등을 종합적으로 평가할 수 있도록 모형을 학습시킬 수 있다는 점이다. 예를 들어, 사전에 정의된 긍정 또는 부정적 단어 목록을 기계적으로 사용하는 경우⁰³ 해당 단어가 사용된 전후 문맥을 고려하지 못할 위험이 있다. 또한, 모형의 개발 단계에서부터 인간 연구진이 개입하고 결과를 검토하는 절차⁰⁴를 통하여 모형의 추론 결과물에서 발생할 수 있는 편의(bias)와 오류(false positive 또는 false negative)를 최소화할 수 있다.

뉴스 데이터를 분류하는 기준은 한국은행 뉴스 심리지수를 참고하여 문장 단위로 긍정, 부정, 중립의 의미를 부여하는 방법을 사용하였다.⁰⁵ 이는 하나의 완성된 글이 있을 때, 내재된 의미를 판단할 수 있는 최소 단위를 문장으로 보고 각 문장이 시장 심리에 대하여 의미하는 바를 크게 세 가지(긍정, 부정, 중립) 기준으로 나눈 것을 의미한다. 라벨링에 사용된 문장 표본은 약 3,000개로 Kim et al.(2024)과 같으며, 분류 기준별 예시 문장은 <표 II -1>과 같다.

<표 II-1> 데이터 라벨링 예시

예시 문장	분류
“이를 반영해 증권가는 연이어 목표가를 하향하고 있다.”	부정
“이어 내년 순이익은 올해보다 33% 추가로 증가할 것으로 예측했다.”	긍정

주 : 시장 동향에 대한 긍정 또는 부정의 예측이 명시적으로 드러나지 않는 문장은 모두 중립으로 분류
 자료: 네이버 뉴스

03 예를 들어, ‘주가의 상승’이라는 단어가 들어간 문장을 기계적으로 긍정적인 투자 심리를 나타내는 문장이라고 분류하는 규칙기반(rule-based) 모형을 들 수 있다.

04 이와 같은 원칙을 흔히 human-in-the-loop(HITL)이라고 한다.

05 한국은행 뉴스 심리지수의 학습 데이터 라벨링 방법론은 서범석 · 이영환 · 조형배(2022)에서 자세하게 논의하고 있다.

문장의 분류 기준은 투자자의 관점에서 뉴스에 명시된 사실이 시장 전반 또는 개별 종목에 호재인 경우를 ‘긍정’, 악재인 경우를 ‘부정’, 판단이 어려운 경우를 ‘중립’으로 하였으며 이는 Kim et al.(2024)의 기준과 일치한다. 다만, 인간 연구진의 분류는 개인의 인식에 따라서 편이가 발생할 수 있는데(서범석 · 이영환 · 조형배, 2022) 이를 최소화하기 위하여 학습 과정에서 ‘부정’과 ‘중립’은 같은 결과로 처리하여 명백하게 ‘긍정’인 경우를 모형이 식별할 수 있도록 학습시키는데 중점을 두었다.

다. 전처리 과정

수집된 뉴스 기사 텍스트 자료는 전처리를 거치게 되는데 이는 여러 단계로 구성된다. 첫 번째 절차는 긴 문서를 문장 단위로 분절하는 과정이다. 일반적인 트랜스포머 기반 LLM은 학습에 사용되는 입력 텍스트의 최대 길이를 제한하고 있어 긴 문서를 문장 단위로 분해할 필요가 있다. 문장을 분해하는 방식으로는 정규식(regular expression)으로 작성된 재량적인 규칙에 기반한 방식과 확률모형에 기반한 방식이 있는데 본 연구에서는 후자를 선택하였다.⁰⁶

<표 II-2> 텍스트 전처리 과정 예시

단계	예시
원문	이날 코스피 지수는 7거래일 연속 상승해 2,671로 거래를 마쳤다.
	↓
정제	이날 코스피 지수는 거래일 연속 상승해 로 거래를 마쳤다
	↓
형태소 분해	[이날, 코스피, 지수, ##는, 거래일, 연속, 상승, ##해, 로, 거래, ##를, 마쳤, ##다]
	↓
단어 임베딩	$C_i = \begin{bmatrix} \text{이날:} & [.453, .830, .865, \dots, .736] \\ \text{코스피:} & [.524, .922, .301, \dots, .528] \\ & \dots \\ \text{##다:} & [.662, .210, .086, \dots, .034] \\ \text{[PAD]:} & [.205, .151, .704, \dots, .408] \end{bmatrix}$

주 : 단어 임베딩에 제시된 수치는 개념의 설명을 위한 가공의 숫자로서 실제와 같지 않을 수 있음
 자료: 네이버 뉴스, 저자 계산

06 문장을 분리하기 위하여 파이선(Python) 모듈인 KSS(Korean String processing Suite)를 사용하였다. 해당 모듈에 대한 설명은 <https://github.com/hyunwoongko/kss>에서 확인할 수 있다.

학습 데이터를 문장 단위로 나눈 후에는 각 문장에 대하여 <표 II -2>와 같은 전처리 과정을 적용한다. 가장 먼저 문장에서 수치 및 특수기호를 제거하고 문자 정보만을 추출하는 정제(cleaning) 단계를 거친다. 이 과정은 다음 단계에서 사용되는 형태소 또는 토큰(token)의 가짓수를 줄여 모형이 뉴스에 포함된 언어 정보의 학습에 집중하도록 한다. 정제된 문장은 최소 구성 단위인 형태소로 분해하게 되는데, 특히 교착어인 한글의 특성에 따라 명사, 동사, 형용사 뿐만 아니라 접사, 조사도 독립된 형태소로 분리된다. 이 때, 문장의 길이가 짧은 경우는 [PAD]를 추가하고 긴 경우는 삭제하여 일관된 길이로 맞춘다.⁰⁷ 마지막으로 각 형태소를 인접 단어와의 연관성을 반영한 수치 벡터로 변환하는 과정을 임베딩(embedding)이라 한다. 이와 같은 과정을 거치면 각 문장을 수치 벡터의 누적인 텐서(tensor)로 표현할 수 있다.

라. 모형의 학습 알고리즘과 결과의 검증

이와 같은 임베딩 과정을 거친 학습 데이터의 집합을 C 라고 하면 모형의 학습은 다음과 같은 함수 관계를 찾는 것으로 표현할 수 있다.

$$f_{\theta} : C \rightarrow P$$

이 때, P 는 유한한 원소를 가지는 레이블(label)의 집합을 의미한다. 구체적으로는 이진(binary) 분류체계를 사용하여 다음과 같이 정의한다.

$$P = \{p^{(0)}, p^{(1)}\}$$

여기서 $p^{(0)}$ 는 부정적, $p^{(1)}$ 은 긍정적인 신호를 의미한다. 따라서 f_{θ} 는 (임베딩 과정을 거친) 문장을 긍정 또는 부정으로 분류하는 함수라고 할 수 있다. 트랜스포머 구조는 이를 매개변수(parameter) θ 를 가지는 중첩된 신경망(neural network) 모형으로 구축한다.⁰⁸ 모형의 학습은

07 토큰의 최대 길이는 128로 설정하였다.

08 임베딩 및 문장 분류 모형은 한국어 분류 및 평가에 최적화된 대형 기반 모형인 KLUE(Korean Language Understanding and Evaluation)를 사용하였다. 해당 모형의 구조와 성능에 대한 논의는 Park et al.(2021)에서 확인할 수 있다.

전체 데이터에서 무작위로 추출하여 긍정 또는 부정의 의미로 사전에 라벨링된 표본⁰⁹에 대하여 다음과 같은 교차 엔트로피(cross entropy) 손실함수(loss function)를 최소화하는 매개변수를 찾는 과정이다.

$$L(\theta|C_{labeled}) = - \sum_{c_i \in C_{labeled}} p_i^{(1)} \log(f_\theta(c_i)) + p_i^{(0)} \log(1 - f_\theta(c_i))$$

학습된 모형 f_θ 가 있을 때, 주어진 문장 c_i 에 대하여 시장 동향에 긍정 또는 부정의 의미를 내포하고 있을 확률은 각각 $f_\theta(c_i), 1 - f_\theta(c_i)$ 로 계산할 수 있다.

<표 II-3> 모형의 확률값 계산 예시

	예시 문장	$f_\theta(c_i)$
1	이날 코스피 지수는 거래일 연속 상승해 로 거래를 마쳤다	0.9812
2	코스피 상승세가 주춤할 것으로 내다보고 굼버스 상품을 대거 사들인 개인투자자의 손실이 누적되고 있다	0.0261

주 : 주어진 예시 문장은 정제를 거쳐서 수치와 특수문자 정보를 제거한 형태임
 자료: 네이버 뉴스, 저자 계산

모형이 라벨링된 데이터를 학습한 이후에는 라벨링이 되어있지 않은 모형에 대해서도 ‘긍정’문 일 가능성을 $f_\theta(c_i)$ 로 수치화할 수 있게 된다. <표 II-3>은 동일한 날짜에 간행된 서로 다른 기사에서 임의로 선정한 두 가지 예시 문장에 대하여 모형이 계산한 확률값을 보여주고 있다. 첫 번째 문장의 경우 대체로 상승장에 대한 기대감이 반영된 문장으로 볼 수 있는데, 모형은 해당 문장에 대한 확률값을 0.9812로 계산하여 ‘긍정’의 신호가 강하다는 점을 정확히 포착하고 있다. 반대로 두 번째 문장은 동일한 시장 환경에서도 ‘개인투자자의 손실’과 같은 부정적인 의미를 내포하고 있어 확률값이 비교적 낮은 0.0261로 계산된 것을 확인할 수 있다.

09 라벨링된 텍스트는 Kim et al.(2024)에서 사용된 표본을 기반으로 하였다.

2. 학습 결과의 해석

자연어에 기반한 초거대 AI 모형은 매우 복잡한 구조를 가지고 있어 입력된 정보와 출력된 결과 값을 명확하게 연결하여 설명하기 어려운 경우가 많다(Danilevsky et al., 2020). 따라서 본 연구에서는 문장별로 계산된 확률값을 각 문장을 구성하는 형태소 단위로 분해하여 해석하는 방안을 제시하여 심리지수 구축 과정의 투명성을 제고한다.

구체적으로 본 연구에서는 Lundeberg & Lee(2007)가 제안한 샤플리(Shapley) 값을 사용하여 개별 형태소의 기여도를 평가한다. 샤플리 값은 복잡한 머신러닝 모형의 결과물을 해석하기 위하여 제안된 개념으로 입력 요소별로 기여도를 산출하고 이들의 합이 항상 결과물과 일치하여 정량화된 기여도를 직관적으로 표현한다는 장점을 가지고 있다. 이에 크고 복잡한 LLM의 설명 가능성을 제고하기 위한 방안 중의 하나로 샤플리 값을 활용하는 시도가 활발하게 이루어지고 있다(Kokalj et al., 2021; Fantozzi & Naldi, 2024). 샤플리 값의 특성과 트랜스포머 기반 LLM에 적용하는 경우 해석 방법은 부록에서 자세하게 논의하였다.

주어진 문장 c_i 와 그에 상응하는 확률값 $f_i = f_\theta(c_i)$ 가 있을 때, 샤플리 값은 이를 문장에 포함된 형태소별 한계 기여도로 분해한다. 구체적으로, 문장 c_i 에 포함된 모든 형태소의 집합을 T_i 라고 할 경우 각 $\tau \in T_i$ 에 대하여 샤플리 값 $\varphi_\tau(f_i)$ 은 다음과 같은 식을 만족한다.¹⁰

$$f_i = \varphi_0 + \sum_{\tau \in T_i} \varphi_\tau(f_i)$$

여기서 φ_0 는 기준값(baseline)으로 평가에 사용된 모든 문장에 대하여 동일한 값을 가진다. 이를 제외하면 샤플리 값은 특정 문장에서 모형이 추정된 긍정문일 확률값 지표에 각 형태소가 기여한 정도를 선형 분해한 것으로 볼 수 있다. 여기서 개별 형태소 τ 에 대하여 φ_τ 는 양 또는 음의 값을 모두 가질 수 있으며, $\varphi_\tau > 0 (< 0)$ 인 경우 τ 는 문장이 긍정(부정)문일 가능성을 높이는 것으로 해석된다.

10 해당 결과는 샤플리 값의 효율성(eficiency) 공리로부터 도출할 수 있다. 샤플리 값의 일반적인 특성은 부록을 참고하기 바란다.

한편, 같은 단어도 문맥에 따라 다른 의미를 가지게 되어 반대 방향의 기여도를 가질 수도 있는데 이는 <그림 II-4>의 예시에서 확인할 수 있다. 우선, 여기서 보이는 예시문은 부정적인 견해를 나타내고 있어 샵플리 값의 총합이 음수로 나타난다. 더불어 앞서 <그림 II-3>의 예시와 유사한 토큰인 ‘상승세’가 이 문장에서는 약하게 부정적인 의미로 해석되는데(-2%), 이후에 등장하는 단어인 ‘주춤’, ‘손실’과 같은 단어와 문맥적인 연관관계가 있기 때문이라고 해석할 수 있다.

이와 같은 예시는 본 연구에서 구축된 모형이 시장 상황을 묘사하는 국문 뉴스 기사의 문맥적 의미를 비교적 정확하게 인지할 수 있으며 이를 바탕으로 심리적인 요인을 효과적으로 정량화할 수 있음을 보여주고 있다.

III. 자본시장 심리지수의 계산 및 특징

1. 일별 심리지수의 구축

이와 같이 문장 단위로 계산된 확률값은 날짜별로 통합하여 다음과 같은 산식으로 일별 지수를 구축하는데 사용된다. 우선, 문장별로 계산된 확률값은 아래의 산식을 적용하여 완전한 긍정문을 1, 부정문을 -1로 표준화한다.¹²

$$x_i = (f_{\theta}(c_i) - 0.5) \times 2$$

다음으로 일별 지수값은 해당일에 공개된 기사에 포함된 모든 문장에 대한 표준화 지수의 평균으로 정의한다. 이를 구체적인 수식으로 표현하면, t 일에 해당하는 모든 문장의 집합을 C_t 라고 할 때, 평균 심리지수는 다음과 같다.

$$CMSI_t^{(mean)} = \frac{1}{|C_t|} \sum_{c_i \in C_t} x_i$$

12 표준화 방식은 최소, 최대값을 기준으로 변경할 수 있으며(예: 긍정을 100, 부정을 0) 본 연구에서는 직관적인 해석을 위해서 완벽한 긍정/부정문이 각각 1/-1을 갖도록 하였다.

한편, 일중 공개된 기사에서 긍정 및 부정적 보도의 분산을 기준으로 지표를 산출할 수도 있으며 이는 다음과 같이 정의된다.

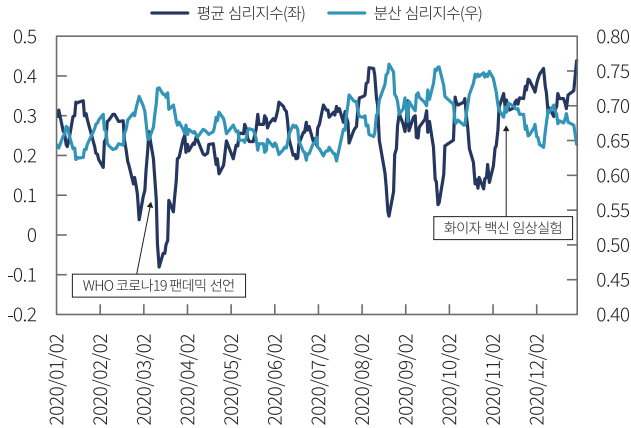
$$CMSI_t^{(std)} = \frac{1}{|C_t|} \sum_{c_i \in C_t} (x_i - CMSI_t^{(mean)})^2$$

평균 심리지수와 분산 심리지수는 서로 다른 측면에서 시장에 대한 심리를 정량화하여 보여준다. 전자의 경우 해당일에 시장의 동향을 묘사하는 뉴스의 평균적인 견해를 수치화한 정보라고 해석할 수 있다. 즉, 평균적으로 투자자 심리 또는 시장 환경에 대하여 긍정적인 묘사가 많은(적은) 날에는 평균 심리지수가 높게(낮게) 나타날 것이다. 반면, 후자의 경우 시장 상황에 대한 견해의 불확실성 정도를 의미하는 것으로 볼 수 있다. 긍정적인 묘사와 부정적인 묘사가 혼재하는 날에는 평균 심리지수가 0에 가깝더라도 분산 심리지수가 높아 해석의 불확실성이 높을 수 있다. 이처럼 다층적인 지표를 구성할 경우 시장 심리에 대한 다면적인 해석이 가능한데 이에 대해서는 III-3.절에서 보다 자세히 설명한다.

2. 시장 심리 동향 분석

본 절에서는 앞서 소개한 방식과 같이 구축된 자본시장 심리지수의 시계열적인 동향을 평균 심리지수를 중심으로 소개하고 의미를 해석한다. 아래 <그림 III-1>은 2020년 한 해 동안 일별 평균 및 분산 심리지수의 변화를 보여주고 있다.

<그림 III-1> 2020년 자본시장 심리지수의 동향



주 : 2020년 한 해 동안 평균 심리지수와 분산 심리지수의
일별 동향을 1주일 단위 이동평균으로 나타내었다.

자료: 저자 계산

2020년 초는 코로나-19 팬데믹의 발생으로 투자 심리가 극단적으로 부정적인 경향을 보인 것을 확인할 수 있다. 국제보건기구(WHO)가 팬데믹을 선언한 2020년 3월 11일에 평균 심리지수는 전일 대비 크게 하락하여 이틀 후인 3월 13일 연중 최저치를 기록하였다. 더불어 같은 기간에 분산 심리지수는 크게 상승하는 것을 <그림 III-1>에서 확인할 수 있다. 다만, 평균 심리지수의 하락 추세는 지속성이 크지 않았는데 이는 각국에서 발표한 경기 부양책의 효과가 심리적인 측면에서 빠르게 반영되었다는 점을 시사하는 것으로 해석할 수 있다. 2020년 하반기에는 일시적인 평균 심리지수의 하락 및 분산 심리지수의 상승이 관찰되었으며 화이자(Pfizer) 백신의 임상실험 성공이 발표된 2020년 11월 9일을 기점으로 평균 심리지수의 상승과 분산 심리지수의 하락이 두드러지게 나타나는 것 또한 확인할 수 있다.

3. 지수의 통계적 특성

가. 지수의 분포와 지속성

이 장에서는 앞서 설명한 방식과 같이 문장별 분류 지표의 평균과 분산에 기반한 일별 심리지수의 통계적 특성을 설명한다.

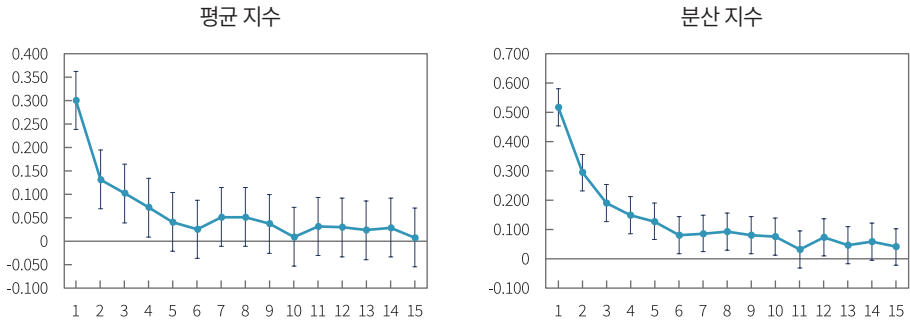
<표 III-1> 심리지수의 기초통계량

	평균	분산	1분위수	중간값	3분위수	왜도	첨도
$CMSI^{(mean)}$	0.215	0.136	0.136	0.233	0.314	-0.731	0.590
$CMSI^{(std)}$	0.725	0.048	0.693	0.728	0.759	-0.183	-0.285

주 : 휴일을 제외한 3,755개의 표본을 대상으로 계산함
 자료: 저자 계산

<표 III-1>은 표본 기간 중 평균 심리지수($CMSI^{(mean)}$)와 분산 심리지수($CMSI^{(std)}$)의 기초통계량을 나타낸다. 우선, $CMSI^{(mean)}$ 의 평균은 0.215로 표본 기간 동안 평균적으로 시장 심리가 긍정적이었음을 알 수 있다. 그러나 해당 지수의 왜도가 음의 값을 가져 심리가 크게 부정적인 날이 있으며 첨도가 0보다 크다는 점에서 이와 같은 극단적인 부정 심리가 나타나는 날의 비중이 적지 않다고 볼 수 있다. 한편, $CMSI^{(std)}$ 의 분포는 평균과 중간값, 분위수의 차이가 적도 첨도가 음의 값을 가지므로 전체 표본 기간 동안 대체로 일관된 값을 가지고 있는 것으로 보인다.

<그림 III-2> 심리지수의 편자기상관함수



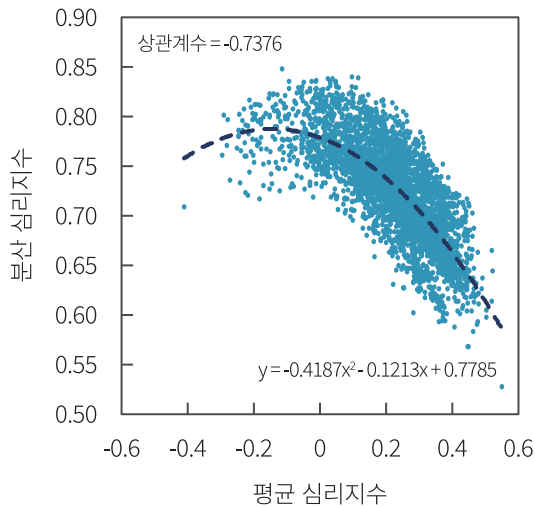
주 : 그림에서 점은 가로축의 시차(lag)에 대응하는 심리지수의 편자기상관계수(partial autocorrelation) 추정값을, 세로선은 95% 유의수준에서 측정한 추정값의 신뢰구간을 의미
 자료: 저자 계산

일별 심리지수는 일정 수준의 지속성(persistence)을 가지는 것으로 나타난다(<그림 III-2> 참고). $CMSI^{(mean)}$ 의 편자기상관함수(partial autocorrelation function)를 보면 최대 4일의 시차까지 양의 상관관계수 값을 보인다. 이는 평균적인 시장 심리가 대체로 1주일 정도 지속적으로 유지됨을 의미한다. 한편 $CMSI^{(std)}$ 는 지속성이 상대적으로 길게 나타나서 약 10일 정도의 시차까지 유의미하게 나타나고 있다. 이와 같은 현상은 시장 환경에 대한 해석의 불확실성이 평균적인 심리의 변화보다 상대적으로 더 오래 지속됨을 의미한다.

나. 평균 및 분산 지수의 관계

평균과 분산 심리지수의 상관관계를 보면, 대체로 반대 방향으로 움직이는 것을 확인할 수 있다 (<그림 III-3> 참고). 전체 표본 기간을 대상으로 일별 $CMSI^{(mean)}$ 와 $CMSI^{(std)}$ 의 분포를 보면 음의 상관관계(상관계수 -0.7376)가 있음을 알 수 있다. 이를 해석하면 시장 심리가 평균적으로 긍정적인 시기에는 일 중 뉴스에 나타난 견해가 대체로 일치하는 반면 평균적으로 부정적인 심리가 나타나는 날에는 상황을 묘사하는 방식의 차이가 커서 시장 환경에 대한 불확실성이 높게 나타나는 것으로 볼 수 있다.

<그림 III-3> 평균 및 분산 지수의 관계



주 : 표본 기간 중 일별 평균 및 분산 심리지수의 분포,
 점선은 이차함수로 근사한 관계를 나타냄
 자료: 저자 계산

다만, 평균과 분산 심리지수의 관계는 선형보다 비선형에 가까운 양상을 보여 해석에 주의할 필요가 있다. <그림 III-3>에서 점선은 두 심리지수의 일별 분포 사이의 관계를 이차함수로 근사한 것인데 평균 심리지수가 0에 가까운 값에서 분산 심리지수가 극대화되는 것을 확인할 수 있다. 이는 일중 뉴스 기사가 전달하는 심리가 긍정과 부정 양면에 고르게 분포하고 있을 때 해석의 불확실성이 높다는 점을 시사한다. 이와 같은 경향을 고려할 때, 평균 심리지수가 중립에 가까운 날은 지수의 해석에 주의를 기울일 필요가 있으며 나아가 지수를 의견의 분산 정도로 표준화한 지표 활용하는 방안을 고려할 수 있다.

IV. 자본시장 심리지수의 의미와 활용 방안

본 연구에서는 뉴스 데이터를 학습한 대형언어모형에 기반하여 주식시장 심리지수를 구축하는 방안을 소개하고 지수의 정보적 특성을 분석하였다. 트랜스포머(transformer) 구조의 범용 모형을 기반으로 약 600만 건의 증권 뉴스 텍스트를 학습하여 시장 참여자의 인식을 정량화한 심리지수를 구축하는 방안을 소개하였다. 특히, 지수의 구축에 필요한 텍스트 전처리 과정과 텍스트 분류 모형의 구조를 상세히 설명하고 문장별로 계산된 확률값에 대하여 샤플리 값 분석을 통해 각 단어의 기여도를 정량적으로 평가하는 방안을 제시하였다. 이와 같은 분석 절차를 통해 지수 산출 과정의 투명성과 해석 가능성을 제고하였다.

문장별로 측정된 긍정 또는 부정 심리에 대한 분류 지표를 표준화하여 평균 심리지수와 분산 심리지수를 일 단위로 구축하였다. 두 지표는 서로 보완적인 정보를 제공하는데, 전자는 시장 심리에 대한 낙관 또는 비관적인 견해의 비중을 의미하는 반면 후자는 일중 견해가 얼마나 분산되었는지를 측정하여 불확실성의 정도를 정량화한 것으로 볼 수 있다.

체계적이고 정량화된 자본시장 심리지수는 근거기반(evidence-based) 의사결정을 보조할 수 있다. 우선, 시장 참여자의 비정형화된 심리적 요인의 동향을 실시간 또는 고빈도로 측정할 수 있어 적시성 면에서 전통적인 지표와 비교하여 우수한 성능을 보인다. 더불어 거시지표 또는 재무 정보만으로는 파악하기 어려운 시장 심리의 변화(과열, 공포, 낙관, 불안 등)를 조기에 감지할 수 있어 보완재로서 유의성이 높다(장보성, 2026). 나아가 시장 참여자의 심리를 정량화하여 주식시장 전반의 정보 효율성과 역동성을 파악하는데 유용할 수 있다(김민기, 2026). 종합적으로, 자본시장 심리지수는 투자 및 정책 의사결정을 지원하는 정보의 통로로서 활용성이 높다고 평가할 수 있다.

끝으로 본 연구에서 소개한 자본시장 심리지수의 고도화를 위한 후속 연구의 방향을 제시하고자 한다. 우선, 다차원 라벨링 체계의 도입을 고려할 수 있다. 본 연구에서는 이원(binary) 분류를 활용하였으나, 향후 뉴스의 영향을 세분화된 범주(예: 매우 긍정, 보통, 중립, 보통 부정, 매우 부정)로 구분하거나 주제별 분류와 결합하면 문맥을 보다 세밀하게 분석하여 정보량이 더욱 풍부해질 것이다. 더불어 학습 알고리즘에서 배제된 정량적 수치 정보를 포괄적으로 인식하는 멀티모달(multi-modal) 모형을 고려할 수 있다. 예를 들어, 노성호 · 이상호(2025)는 정량적 공시인

재무제표와 정성적 정보인 주식 공시문을 동시에 학습시킨 언어모형을 통해 부실 징후를 조기에 식별할 수 있음을 제시하였는데 이를 뉴스 텍스트의 학습에도 적용할 수 있다. 마지막으로 뉴스의 일 중 시간대별 분포를 고려하여 지수를 세분화하여 보다 고빈도의 지표를 구축할 수 있다. 이와 같은 후속 연구를 통해 자본시장 심리지수의 활용성을 높이고 정보의 효율적인 확산에 실질적으로 기여할 수 있을 것으로 기대한다.

참고문헌

김민기, 2026, 『자본시장 심리지수 시리즈 3: 투자자 심리와 주식시장의 관계 고찰』, 자본시장 연구원 이슈보고서 26-03.

김현중 · 임중호 · 이해영 · 이상호, 2019, 온라인 뉴스 기사를 활용한 경제심리보조지수 개발, 『국민계정리뷰』 2019(2), 1-33.

노성호, 2024, 『증권업 경쟁력 강화 시리즈 2: 대형언어모형의 발전과 금융정보분석에의 활용 방안』, 자본시장연구원 이슈보고서 24-02.

노성호 · 이상호, 2025, 『머신러닝을 활용한 재무제표 정보의 유용성 평가: 부실 징후의 조기 탐색을 중심으로』, 자본시장연구원 연구보고서 25-04.

서범석 · 이영환 · 조형배, 2022, 기계학습을 이용한 뉴스심리지수(NSI)의 작성과 활용, 『국민계정리뷰』 2022(1), 68-90.

장보성, 2026, 『자본시장 심리지수 시리즈 2: 거시 · 금융변수와의 관계와 유용성』, 자본시장 연구원 이슈보고서 26-02.

네이버 뉴스 <https://news.naver.com/>

Baker, M., & Wurgler, J., 2007, Investor sentiment in the stock market, *Journal of Economic Perspectives* 21(2), 129-151.

Benhabib, J., Liu, X., Wang, P., 2016, Sentiments, financial markets, and macroeconomic fluctuations, *Journal of Financial Economics* 120(2), 420-443.

Bybee, L., Kelly, B., Manela, A., Xiu, D., 2024, Business News and Business Cycles, *The Journal of Finance* 79(5), 3105-3147.

Bybee, L., Kelly, B., Su, Y., 2023, Narrative asset pricing: Interpretable systematic risk factors from news text, *The Review of Financial Studies* 36(12), 4759-4787.

Da, Z., Engelberg, J., Gao, P., 2015, The Sum of All FEARS Investor Sentiment and Asset Prices, *The Review of Financial Studies* 28(1), 1-32.

- Danilevsky, M., Qian, K., Aharonov, R., Katsis, Y., Kawas, B., Sen, P., 2020, A survey of the state of explainable AI for natural language processing, *arXiv preprint arXiv:2010.00711*.
- Du, K., Xing, F., Mao, R., Cambria, E., 2024, Financial sentiment analysis: Techniques and applications, *ACM Computing Surveys* 56(9), 1-42.
- Fantozzi, P., Naldi, M., 2024, The explainability of transformers: Current status and directions, *Computers* 13(4), 92.
- Huang, C., Simpson, S., Ulybina, D., Roitman, A., 2019, News-based sentiment indicators, *IMF Working Paper* 19/273.
- Huang, D., Jiang, F., Tu, J., Zhou, G., 2015, Investor sentiment aligned: A powerful predictor of stock returns, *The Review of Financial Studies* 28(3), 791-837.
- Kim, S., Choi, Y., Jeon, J.H., Lu, Y., 2024, Sentiment Matters in Stock Market: Construction of Sentiment Index Using Machine Learning, *Journal of Economic Theory and Econometrics* 35(4), 87-112.
- Kokalj, E., Škrlj, B., Lavrač, N., Pollak, S., Robnik-Šikonja, M., 2021, BERT meets Shapley: Extending SHAP explanations to transformer-based classifiers, *Proceedings of the EAACL hackashop on news media content analysis and automated report generation*, 16-21.
- Lundberg, S.M., & Lee, S.I., 2017, A unified approach to interpreting model predictions, *Advances in Neural Information Processing Systems* 30.
- Park, S., Moon, J., Kim, S., Cho, W.I., Han, J., Park, J., Song, C., Kim, J., Song, Y., Oh, T., Lee, J., Oh, J., Lyu, S., Jeong, Y., Lee, I., Seo, S., Lee, D., Kim, H., Lee, M., Jang, S., Do, S., Kim, S., Lim, K., Lee, J., Park, K., Shin, J., Kim, S., Park, L., Oh, A., Ha, J.-W., Cho, K., 2021, KLUE: Korean language understanding evaluation, *arXiv preprint arXiv:2105.09680*.
- Shapley, L.S., 1953, A Value for n-Person Games, In Kuhn, H.W., Tucker, A.W. (Eds.), *Contributions to the Theory of Games Vol. II*, 307-318, Princeton: Princeton University Press.

<부록> 샤플리 값의 정의와 특성

샤플리(Shapley) 값은 협력적 게임(cooperative game)에서 참가자들에게 이익 또는 비용을 공정하게 분배하는 방법에 대한 개념으로 제안되었다(Shapley, 1953). 개별 참가자에 대응하는 샤플리 값은 해당 참가자를 포함한 모든 연합(coalition)에서 발생 가능한 결과에 대하여 한계적 기여도(marginal contribution)를 평균한 값이다. 이를 구체적으로 정의하면 다음과 같다. 게임 참가자의 집합을 $T = \{1, 2, \dots, n\}$ 이라 하고 모든 부분집합—즉, 연합(coalition)—인 $S \subseteq T$ 에 대하여 결과를 측정하는 특성함수(characteristic value function)를 $f : 2^T \rightarrow R$ 이라고 할 경우¹³ 참가자 $\tau \in T$ 에 대응하는 샤플리 값은 다음과 같다.

$$\varphi_\tau(f) = \sum_{S \subseteq T \setminus \{\tau\}} \frac{|S|!(n-|S|-1)!}{n!} (f(S \cup \{\tau\}) - f(S))$$

Shapley(1953)는 위와 같이 정의된 분배 규칙이 다음과 같은 네 가지 공리(axiom)를 만족시키는 유일한 규칙임을 증명한 바 있다.

효율성(efficiency): 결과값이 모든 참여자에게 빠짐없이 분배됨을 의미한다.

$$\sum_{\tau \in T} \varphi_\tau(f) = f(T)$$

대칭성(symmetry): 두 참가자가 동일한 역할을 한다면 샤플리 값에 따른 보상 또한 같아야 한다.

$$f(S \cup \{\tau_1\}) = f(S \cup \{\tau_2\}) \Rightarrow \varphi_{\tau_1}(f) = \varphi_{\tau_2}(f)$$

더미(dummy) 또는 널 참가자(null player): 참가자가 게임의 결과에 영향을 미치지 않는다면 샤플리 값은 0이어야 한다.

13 일반적으로 공집합 \emptyset 에 대하여 $f(\emptyset) = 0$ 으로 정의한다.

$$f(S \cup \{\tau\}) = f(S) \Rightarrow \varphi_\tau(f) = 0$$

선형성(linearity): 두 게임의 합에 대한 기여도는 각 게임에 대한 기여도의 합과 같다.

$$\varphi_\tau(af + bg) = a\varphi_\tau(f) + b\varphi_\tau(g)$$

Lundberg & Lee(2017)는 샤플리 값의 개념을 확장하여 머신러닝 또는 AI 모형의 결과값을 입력 데이터의 요소별 기여도로 분해하여 설명하는 방안을 제시하였다. 특히 최근에는 복잡한 LLM의 출력물을 입력 자료를 바탕으로 이해하려는 시도로서 샤플리 값에 기반한 설명가능성에 대한 연구가 활발하게 진행되고 있다(Kokalj et al., 2021; Fantozzi & Naldi, 2024). 샤플리 값을 LLM에 적용할 경우 앞서 설명한 개념을 대응시켜 해석할 수 있다.

<부록 표> LLM 대응 요소

원개념	기호	LLM 대응 요소	설명
참가자	τ	입력 토큰	입력문의 구성 토큰(단어 또는 형태소)를 게임의 참가자로 해석
특성함수	f	확률 예측값	학습된 모형에 기반하여 계산된 긍정/부정 확률값을 게임의 결과로 해석
샤플리 값	$\varphi_\tau(f)$	토큰별 기여도	각 토큰별 한계 기여도로 계산된 확률값을 분해

<부록 표>는 샤플리 값과 관련한 원개념을 LLM의 구성 요소 사이의 대응 관계를 나타낸다. 학습된 LLM을 협력적 게임으로 해석할 경우 입력문을 구성하는 요소인 토큰은 게임의 참가자이고 출력물인 긍정문일 확률은 게임의 결과가 된다. 따라서 샤플리값은 특정한 토큰—경우에 따라서 단어 또는 형태소가 될 수 있다—이 포함되었을 경우 예측된 확률값의 변화를 가능한 모든 토큰의 조합에 대하여 평균하여 계산한 결과를 의미한다. 따라서 $\varphi_\tau > 0$ 인 경우는 긍정일 확률을 증가, $\varphi_\tau < 0$ 인 경우는 반대로 감소시켰음을 알 수 있다. 나아가 주어진 문장에 포함된 모든 토큰에 대한 샤플리 값의 합은 효율성 공리에 따라 문장의 확률 예측값과 동일하게 된다.