

뉴스 미디어 AI 모델링



Schedule Summary

- 4월 말 완료 예정
- 머신러닝 학습데이터 구축
 - 기계학습 DB: 2021년 4월 1일 10만개 구축 완료
 - > ~04.30일까지 샘플링 테스트 및 조정 진행 예정
- 기술 파트 (A.I 모델링)
 - 최종 모델을 가정하여 모델의 정확도 예측을 위한 목표 수립
 - 구축된 80,000개의 데이터 기반으로 데이터 품질 테스트
 - > ~04.30일까지 샘플링 테스트 및 검수 기준 (QC) 확정

(참고) QESG 뉴스 데이터베이스 접근권 공유 완료 및 필요한 소스 코드 제공 완료
- 크롤링 (빅데이터 팀)
 - 구글 뉴스 기반 크롤러 현재 약 20% 정도 진행
 - > 4월 중순 완료 예상
- 플랫폼 / DB 설계 (빅데이터 팀)
 - 프로토타입 설계 위한 세부내역서 확인 및 랜딩 페이지 확정
 - > 4월 말 완료 예상
- 플랫폼의 경우 2단계 5월 부터 추가 개발 진행 예정이며, 연 내에 서비스 정책 확정 및 유료화 론칭 예상

1. 목적 및 기대효과

1. 사람의 업무를 AI로 대체하여 Time cost를 줄이기 위함.
- 사람의 감성 소통을 흉내내는 social intelligence, social analysis를 기계에 학습시킴으로써 Time cost를 획기적으로 줄일 수 있음.

기존 업무, AI 대체 후 업무 COST 비교 (1달 기준, 3개월 치 데이터 처리 가능)

내용	Previous			After applying AI		
		투입 인원	Cost (단위: 천원)		투입 인원	Cost (단위: 천원)
인원 별 cost	인턴		1,600	인턴		1,600
	연구원		4,000	연구원		4,000
작업 내용	뉴스 검색	인턴 4명 (weight=0.2)	$6,400 * 0.2 = 1,280$	뉴스 검색	뉴스미디어 AI Framework	0
	뉴스 수집	인턴 4명 (weight=0.4)	$6,400 * 0.4 = 2,560$	뉴스 수집		0
	뉴스 분류	인턴 4명 (weight=0.4)	$6,400 * 0.4 = 2,560$	뉴스 분류		0
	결과 검수	연구원 1명	1,000	결과 검수	인턴 1명(weight = 0.8) + 연구원 1명(weight = 0.2)	$1,600 * 0.8 + 4,000 * 0.2 = 2,080$
Total cost	<u>7,040</u>			<u>2,080</u>		

1. 목적 및 기대효과

2. 일관성이 없는 사람을 일관성 있는 AI로 대체하기 위함.

- 인공지능은 지치지 않음. 사람에 비해 에러 발생의 빈도나 가능성이 인간에 비해 없음. 사람의 경우 과부하로 인해 끊임없이 일을 할 수 없지만, 인공지능의 경우 생산성 저하, 능률을 걱정할 필요가 없음.
- 반복적인 작업 시 사람의 에러는 일관성 없이 발생하고 이는 조사자의 편향으로 일반적으로 일어남. 하지만 AI의 경우 수치화 된 벡터 값을 학습된 패턴을 통해 인식함으로 AI 에러는 일관성 있게 발생한다. 그러므로 AI의 일관성 있는 에러는 유지보수를 통해 감소 시킬 수 있지만 사람의 에러는 일관성 없이 발생하므로 이를 해결하기 위해 반복적 교육과 많은 시간이 필요한 차이가 있다.

3. Daily ESG rating model – alpha model(risk detecting)

- 현재 반기별로 나오는 대신경제연구소의 ESG rating model에 daily alpha를 추가하는 효과. 다수의 평가기관에서는 수작업으로 정성평가를 통해 데이터를 수집함으로 Daily rating을 할 수 없는 한계점이 있음.
- ESG rating model의 등급이 daily로 바뀔 경우 다양한 파생상품을 생성하여 활용 가능할 것이고 이에 대한 수요가 많을 것으로 예상(반기가 아닌 데일리로 점수를 제공할 경우, 시간가치, 변동성 등을 생성할 수 있음)
- 현재는 대다수의 기관에서 반기, 분기 등으로 rating을 하기 때문에 평가사로 부터 적합판정을 받기 어려운 한계가 있는 것으로 보임. 뿐만 아니라 ESG 기반의 자산이 시장 데이터를 생성할 만큼 많지 않고 국고채, KOSPI 지수와 같은 대표 지수가 없음. 이는 daily로 market data를 생성할 수 없기 때문인 것으로 보임.

e.g. ESG 연동 파생상품(금리스왑)

- 일반적인 스왑 거래에 ESG rating을 연동하여 만기 일시지급 또는 스왑 금리에 가산하여 지급하는 방법 논의 가능.
- 스왑거래 시, Third party KPI 이용 관련한 side letter를 함께 체결해야 하는데, ESG 연동 파생상품은 평가기관의 점수를 이용하여 만기 일시 지급 또는 가산스왑금리 책정 여부를 판단.
- 부적합판정 이유: 반기 별 ESG rating을 파생상품에 연동할 경우 데일리 평가 값을 도출할 수 없기 때문에 스왑을 할 수 없음.
- 데일리 평가 값을 산출하고 이를 연동하여 스왑거래 적합성을 확보 가능.
- 실제로 ESG를 연동한 CDS, IRS 등의 상품 trade의 수요가 있음.

1. 목적 및 기대효과

3. Daily ESG rating model – alpha model(risk detecting)

1. Alpha model scenario 1.

- Rating model 내에 alpha model 삽입.
- 단점: 변동성이 너무 작을 수 있음.

e. g.

E1	E2	E3	S1	S2	S3	G1	G2	G3	G4
0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1

Sum of weight = 1, 즉 1가지

- 붉은색 컬럼들을 Daily로 추적하여 Daily 평가 진행 ^{모델}
- 하얀색 컬럼들은 이전처럼 반기별로 평가 진행

2. Alpha model scenario 2.

- 반기 별 Original ESG rating model + ESG risk alpha model
- 장점: 기존 ESG rating model을 두고 추가 연구를 통해 변동성 및 Turnover 제한하는 model 개발 가능
- 단점: 연구 과정에 있어서 Time cost 소모

E1	S1	S2	G1	G2	+	E1	E2	S1	S2	G1
0.2	0.2	0.2	0.2	0.2		0.2	0.2	0.2	0.2	0.2

Sum of weight = 1,

Sum of weight = 1,

2개 모델

<Original ESG rating model>

<ESG risk alpha model>

2. 국내 경쟁사 비교

	대신경제연구소	지속가능발전소	KCGS
Stewardship code	O	O	O
Scoring framework by using media data	O	O	X
Rating cycle	Daily로 변경 예정, 현재 semiannually	Daily	Annually, 단 수시 조정 있음.
Applying AI	O	O	O
Market overview service	개발예정	O	X
ESG news alert	개발예정	O	X
Number of indicators	206	약 200여개	220
ESG, SRI	O	O	O
ESG performance, incident analysis framework	X	O	X

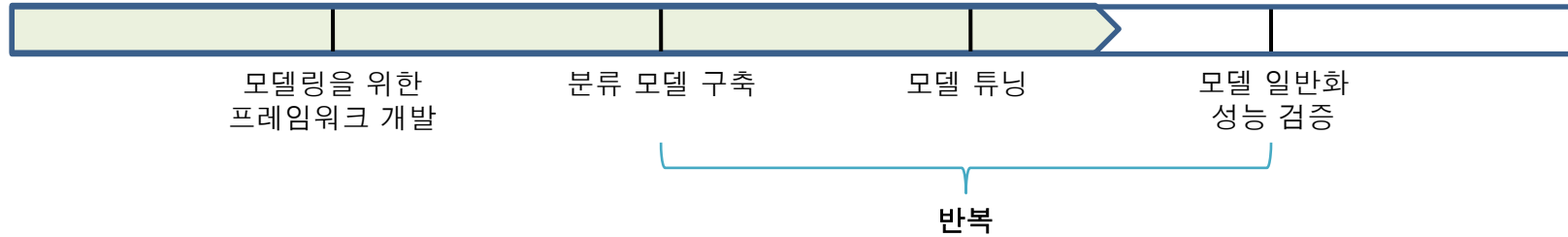
대신경제연구소만의 장점

- 클래스 불균형 문제를 해결하기 위해서 모델을 계층적으로 구성하여 더 세밀한 분류 및 확장이 가능하다. ESG가 아닌 기사를 정교한 알고리즘을 통해 분류 하기 위해서 모든 계층 단계에서 ESG가 아닌 기사를 필터링 하는 구조로 구성되어 있다.
- 지속가능발전소는 일반화된 모델을 사용하여 ESG에 특화된 대신경제연구소의 모델에 비해 정교하지 않음.
- 결론적으로, 지속가능발전소의 특허내용에 의하면 가장 기본적인 워크플로우이며 확장성이 대신경제연구소에 비해 용이하지 않고 다양한 ***ESG risk alpha model***을 추가하는데 있어 어려움이 있을 것으로 예상.

3. 진행 상황

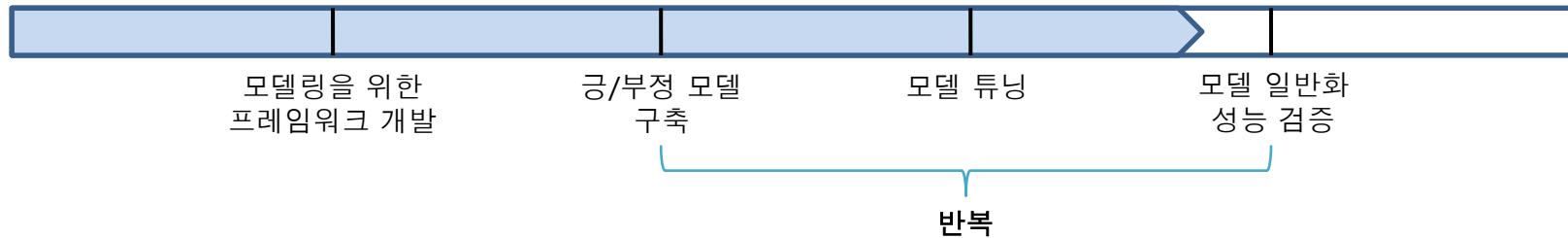
- **분류 모델 : 70% 완성**

정의: 뉴스 미디어 데이터를 Category에 맞게 분류하는 AI 모델



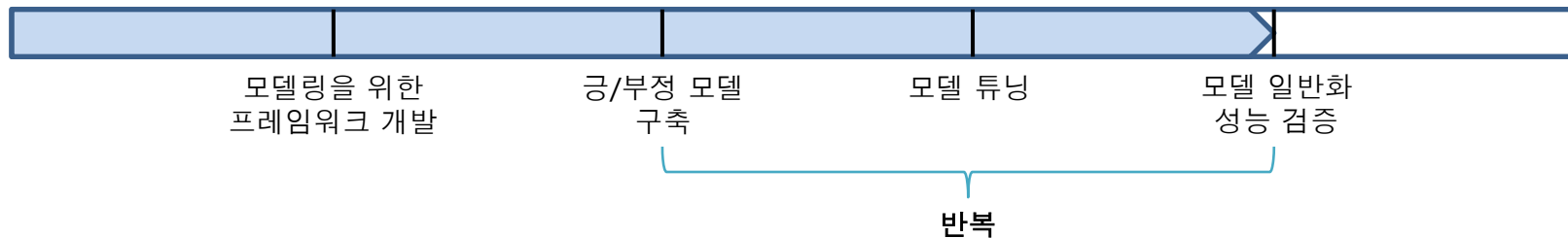
- **긍정/부정 모델: 75% 완성**

정의: 뉴스 미디어 데이터의 긍정/부정을 예측하는 AI 모델



- **기업명 추출기: 80%**

정의: 뉴스 미디어 데이터에서 Coverage 기업이름을 추출하는 AI 모델

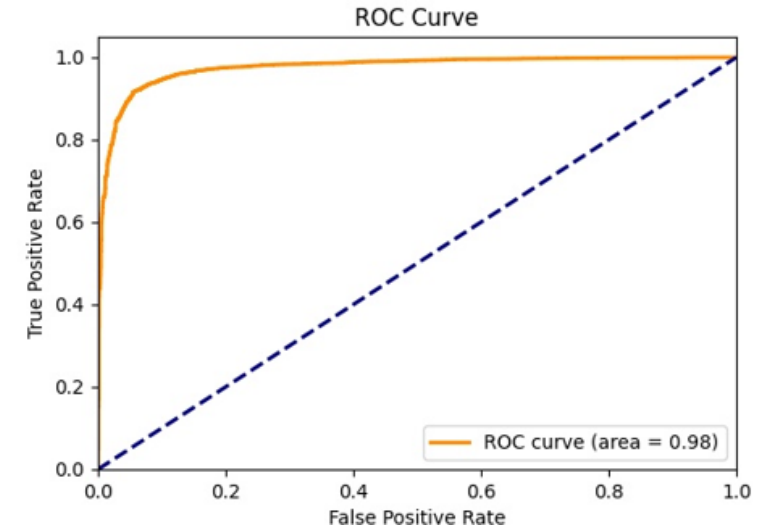


4. 모델 성능 및 산출물

긍정/부정 모델링 결과

Indicator	
f beta	0.960326
AUC	0.975276
F1 score	0.960718
accuracy	0.939449

Confusion Matrix			
		Predicted condition	
		Positive	Negative
Actual condition	Positive	2087 (0.87)	323 (0.13)
	Negative	312 (0.04)	7765 (0.96)

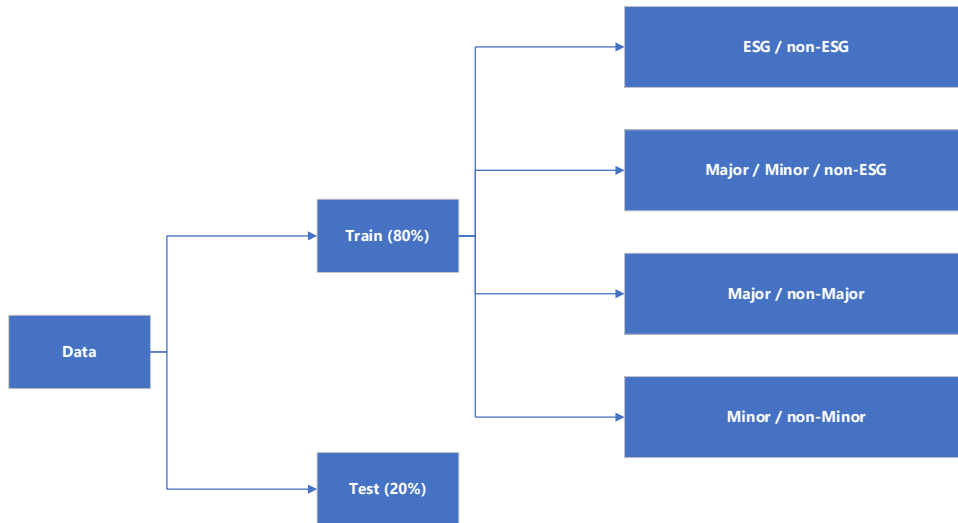


4. 모델 성능 및 산출물 : 카테고리 분류

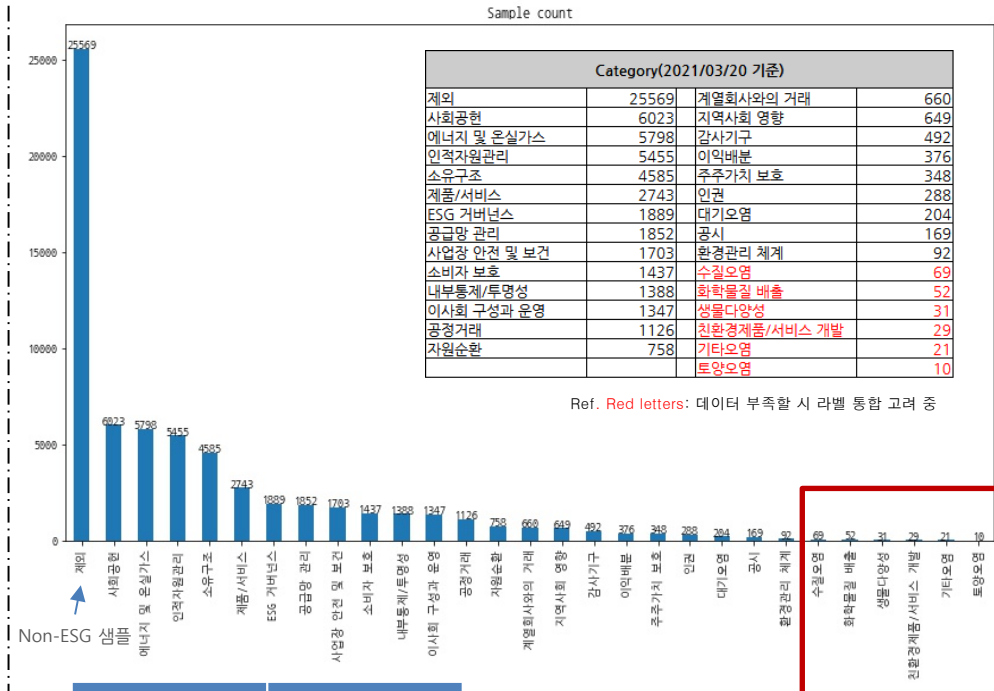
데이터 구성

트레이닝, 테스트 데이터 구성

Holdout 방식으로 모델 평가를 위한 데이터 구성



데이터 클래스 분포



Major	Minor
사회공헌 에너지 및 온실가스 소유구조 제품/서비스 ESG 거버넌스 내부통제/투명성 계열회사와의 거래 감사기구	공급망 관리 사업장 안전 및 보건 소비자 보호 이사회 구성과 운영 공정거래 자원순환 지역사회 영향 이익배분 주주가치 보호 인권 대기오염 공시 환경관리 체계

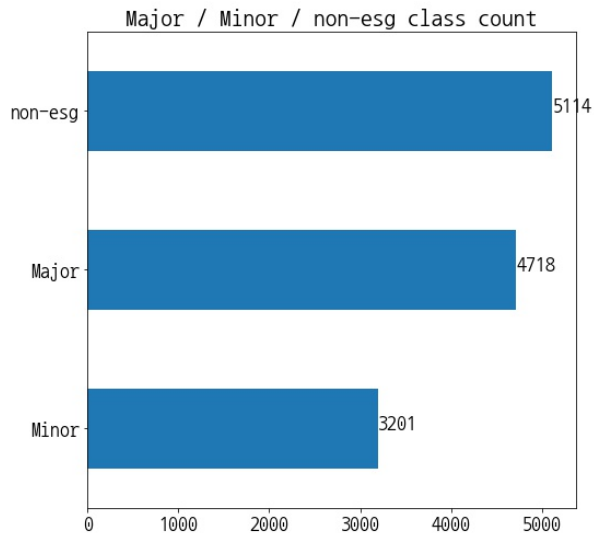
카테고리 구성

샘플 수 문제로 하나로 합침.

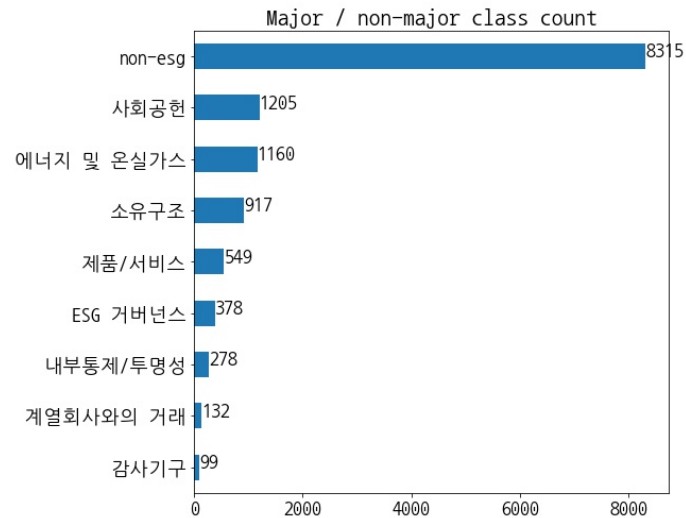
4. 모델 성능 및 산출물 : 카테고리 분류

Task 별 데이터 분포

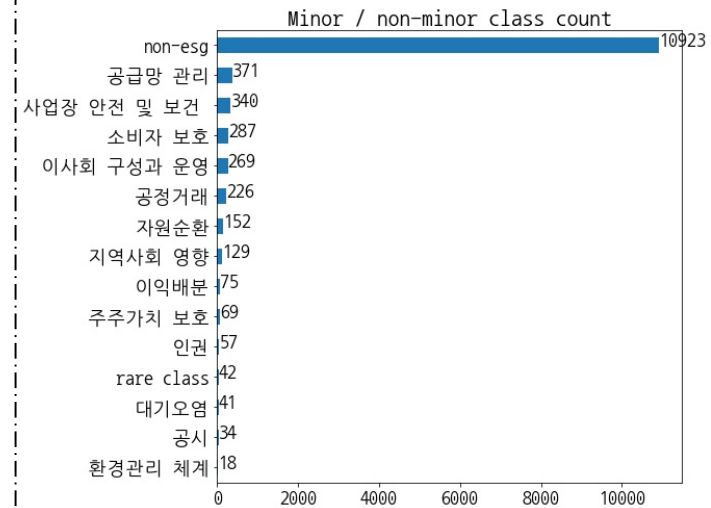
Major / Minor / non-esg : class count



Major / non-Major : class count



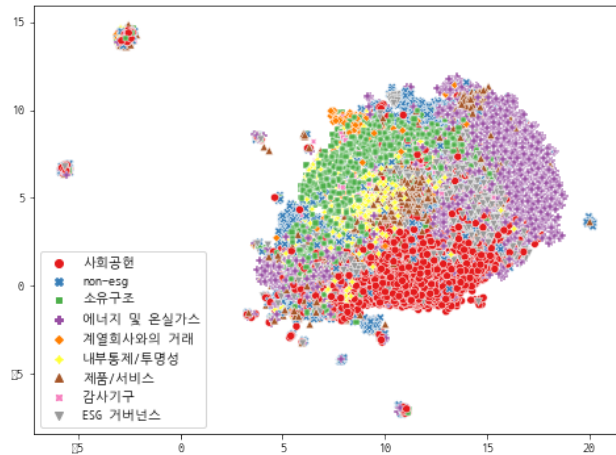
Minor / non-Minor : class count



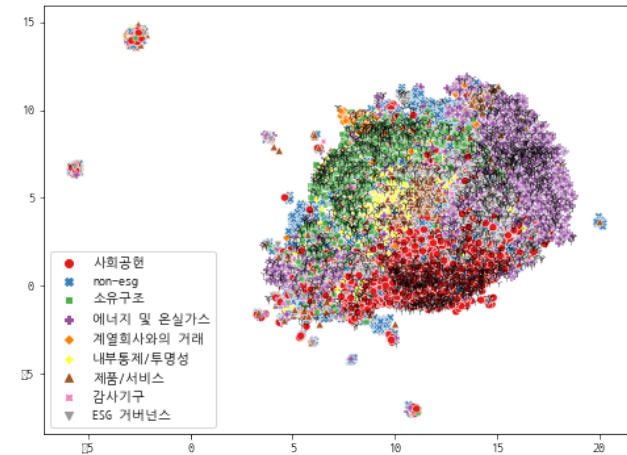
4. 모델 성능 및 산출물 : 카테고리 분류

데이터 품질 확인 : 클러스터링(UMAP with Jaccard Metric)

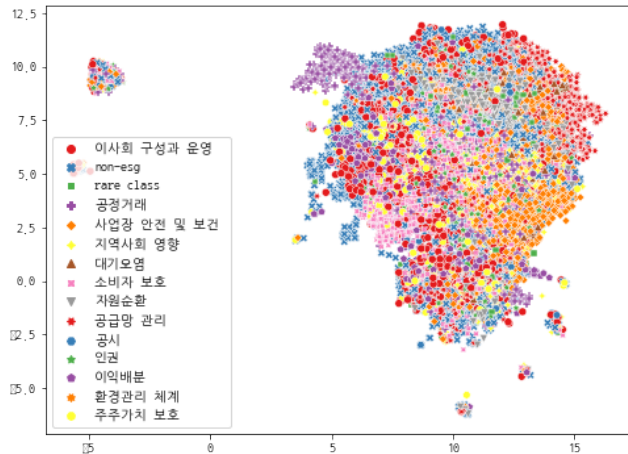
Major Train set only



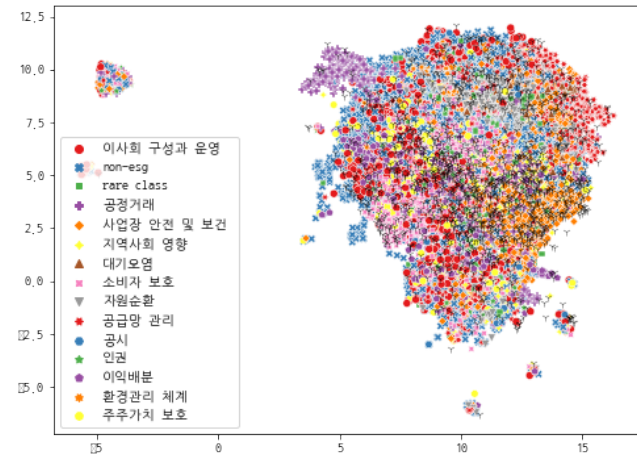
Major Train/Test(black) set only



Minor Train set only



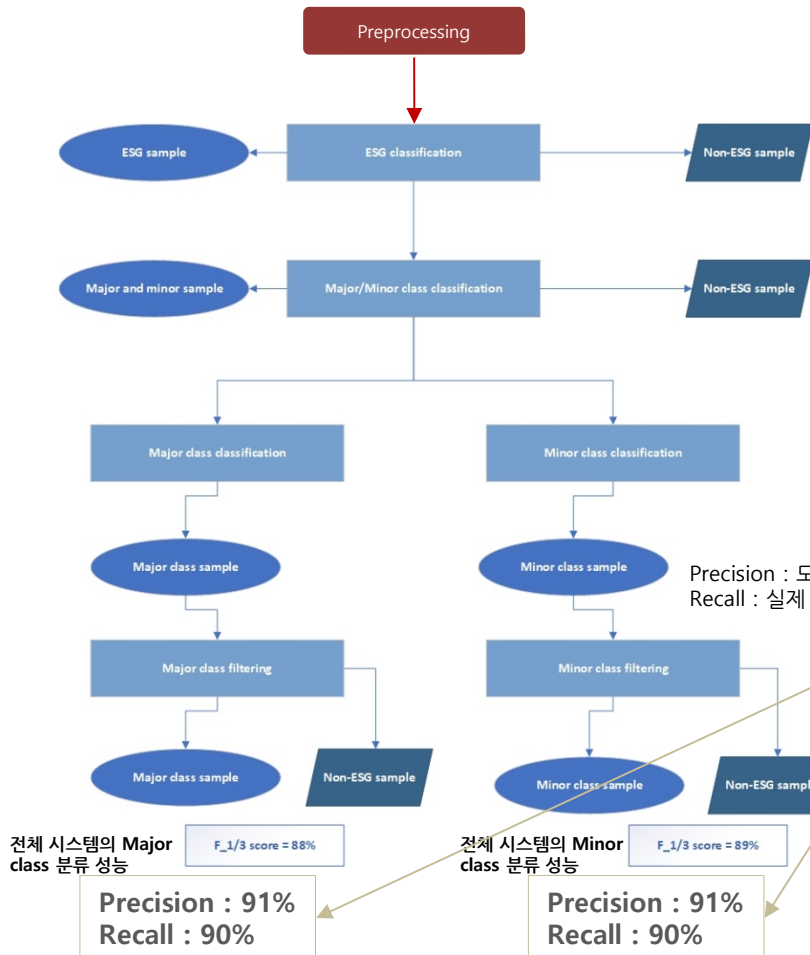
Minor Train/Test(black) set only



4. 모델 성능 및 산출물 : 카테고리 분류

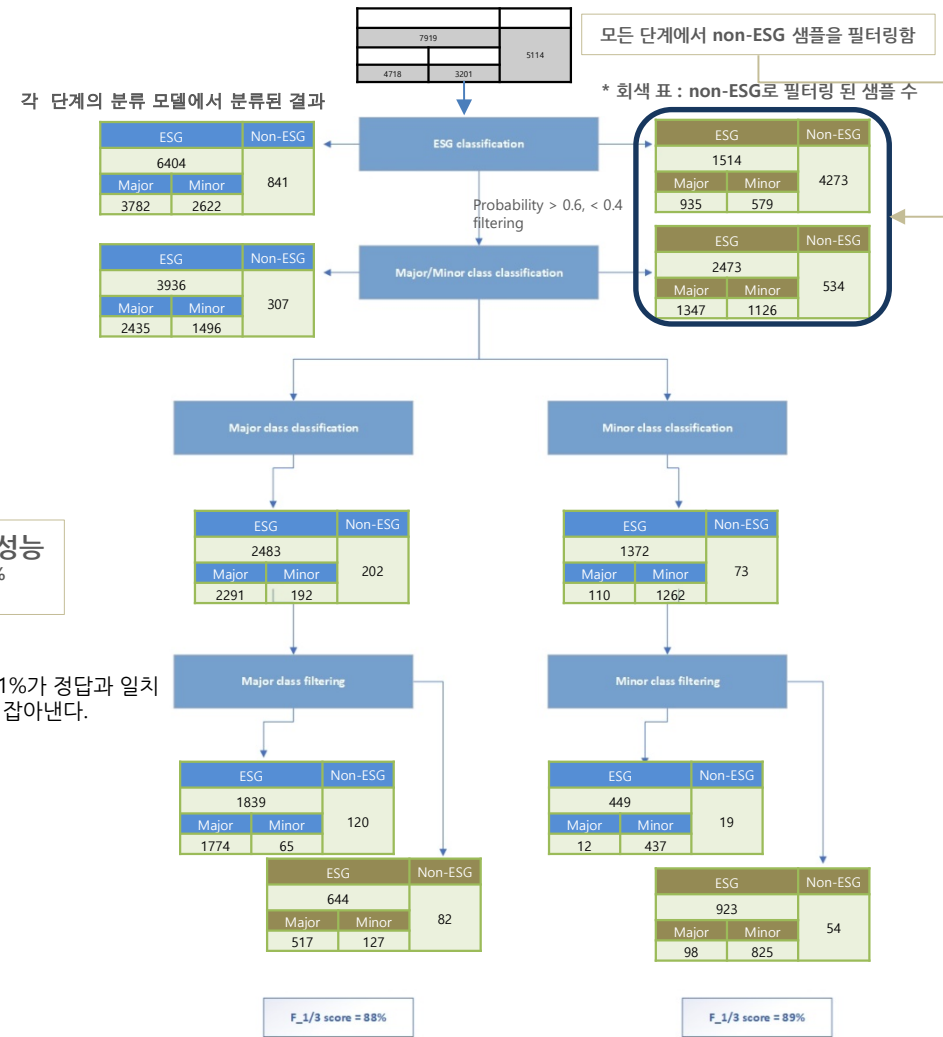
시스템 구성도, 단계별 분류 예시

- 시스템 최종 평가는 f_beta (beta = 1/3) 로 설정
- 평균의 방식은 micro average 사용



ESG 예측 성능
Precision : 95%
Recall : 31%

Precision : 모델이 예측한 카테고리 중 91%가 정답과 일치
Recall : 실제 카테고리 중 90%를 모델이 잡아낸다.



4. 모델 성능 및 산출물 : 카테고리 분류

모델 별 성능

Class name

- 0 제외
- 1 사회공헌
- 2 에너지 및 온실가스
- 3 인적자원관리
- 4 소유구조
- 5 제품/서비스
- 6 ESG 거버넌스
- 7 공급망 관리
- 8 사업장 안전 및 보건
- 9 소비자 보호
- 10 내부통제/투명성
- 11 이사회 구성과 운영
- 12 공정거래
- 13 자원순환
- 14 계열회사와의 거래
- 15 지역사회 영향
- 16 감사기구
- 17 이익배분
- 18 주주차지 보호
- 19 인권
- 20 대기오염
- 21 공시
- 22 환경과 사회 체계
- 23 수질오염
- 24 화학물질 배출
- 25 생물다양성
- 26 친환경제품/서비스 개발
- 27 기타오염
- 28 토양오염

Class 23으로
합침

ESG classification performance

	precision	recall	f1-score	support
0	0.74	0.84	0.78	5114
1	0.88	0.81	0.84	7919
accuracy			0.82	13033
macro avg	0.81	0.82	0.81	13033
weighted avg	0.83	0.82	0.82	13033

Comment
성능이 recall/precision 모두
0.8 이상 균형있게 나온다.

Major/Minor/non-ESG classification performance

	precision	recall	f1-score	support
0	0.77	0.76	0.77	2622
1	0.77	0.86	0.81	3782
2	0.42	0.23	0.29	841
accuracy			0.75	7245
macro avg	0.65	0.62	0.62	7245
weighted avg	0.73	0.75	0.74	7245

Comment
Non-esg 는 성능이 떨어지나, 이후 분류 모델에서도 필터링 하므로,
Major/minor에서의 성능이 중요. 이 둘의 recall/precision 모두
0.75 이상으로 균형 잡힌 성능을 보여준다. .

Major classification performance

	precision	recall	f1-score	support
-1	0.50	0.01	0.02	185
1	0.90	0.99	0.94	612
2	0.87	0.99	0.92	551
4	0.86	0.98	0.92	388
5	0.80	0.84	0.82	38
6	0.89	0.87	0.88	134
10	0.88	0.67	0.76	21
14	0.83	0.65	0.73	23
16	0.50	0.29	0.36	7
accuracy			0.87	1959
macro avg	0.78	0.70	0.71	1959
weighted avg	0.84	0.87	0.83	1959

Comment
Non-esg 는 성능이 떨어진다. 이는 다른
카테고리로 분류되므로 전체적인 성능에
악영향을 미치고있다. 그러나 16 감사기구,
14 계열회사와의 거래 이외의
카테고리에서는 precision/recall모두 0.8
이상으로 균형 잡힌 성능을 보여주고 있다.

Minor classification performance

	precision	recall	f1-score	support
-1	0.92	0.73	0.81	872
7	0.62	0.85	0.71	134
8	0.60	0.91	0.72	164
9	0.55	0.89	0.68	57
11	0.49	0.73	0.59	56
12	0.72	0.60	0.65	30
13	0.81	0.96	0.88	27
15	0.45	0.36	0.40	25
17	0.76	0.80	0.78	20
18	0.80	0.50	0.62	8
19	0.69	0.43	0.53	21
20	0.67	0.57	0.62	7
21	0.44	0.88	0.58	8
22	0.00	0.00	0.00	1
23	0.75	0.20	0.32	15
accuracy			0.75	1445
macro avg	0.62	0.63	0.59	1445
weighted avg	0.80	0.75	0.76	1445

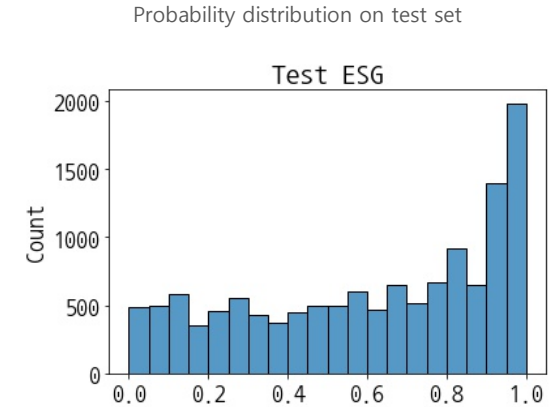
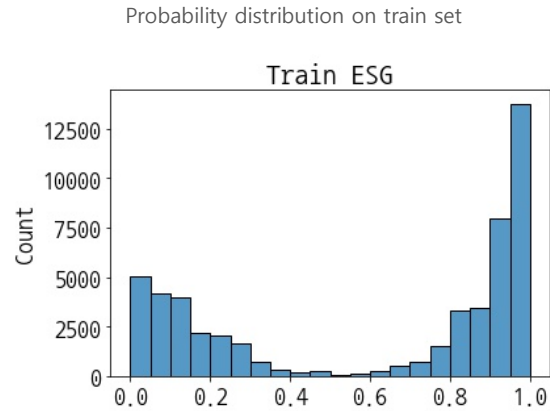
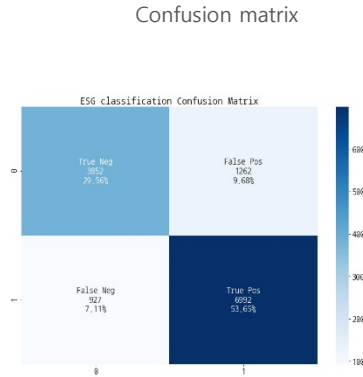
Comment
Non-esg 는 잘 잡아낸다.
카테고리에서는 특별한 패턴을 발견하기
힘든다.

4. 모델 성능 및 산출물 : 카테고리 분류

모델 별 Trustworthy test : confusion matrix, 예측 확률 분포

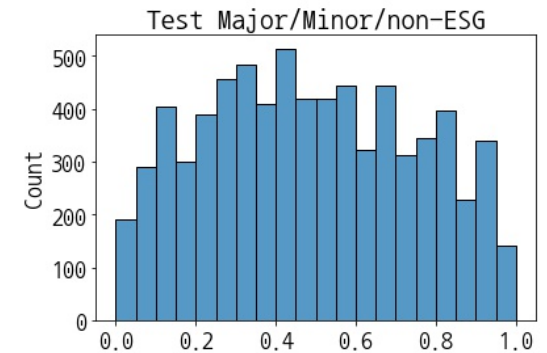
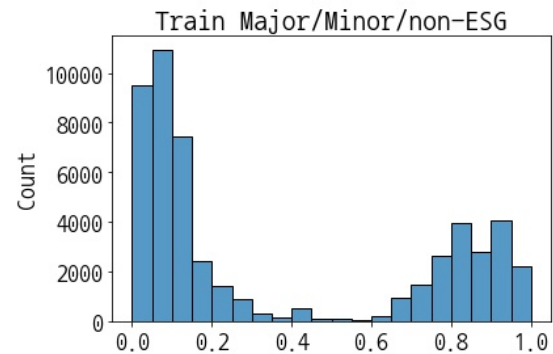
ESG classification performance

Class name
0 : non-ESG
1 : ESG



Major/Minor/non-ESG classification performance

Class name
0 : Minor
1 : Major
2 : non-ESG

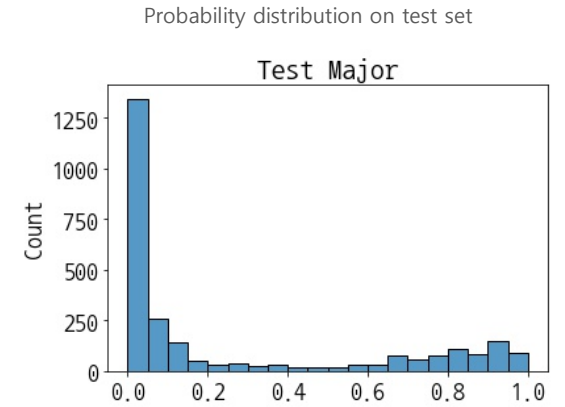
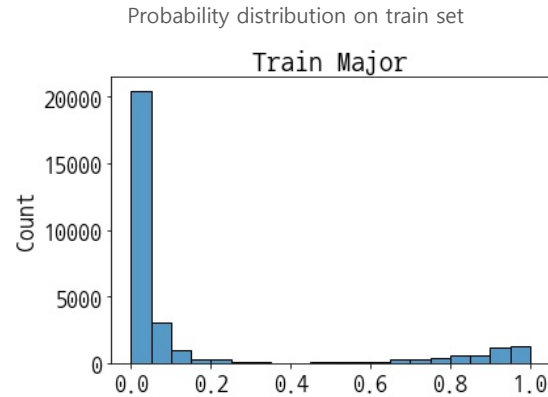
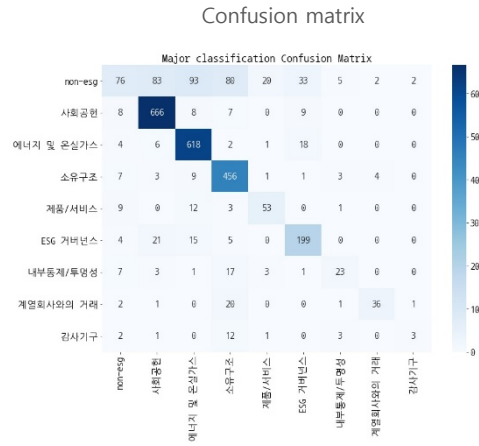


언더 피팅이 의심됨.

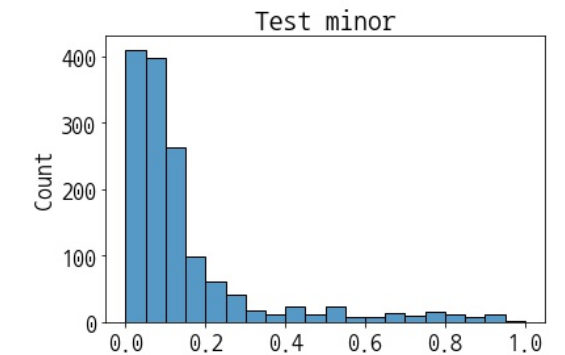
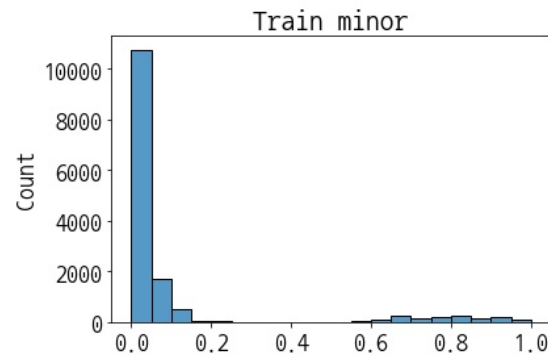
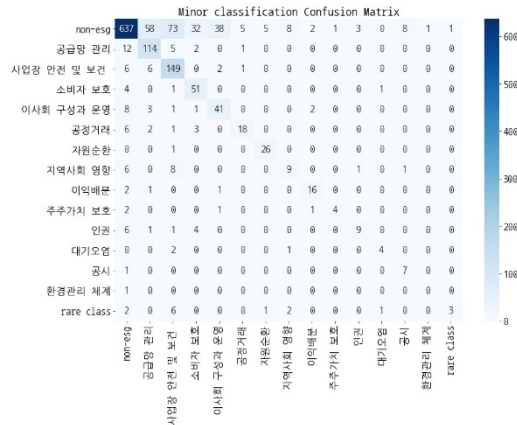
4. 모델 성능 및 산출물 : 카테고리 분류

모듈 별 Trustworthy test : confusion matrix, 예측 확률 분포

Major classification performance



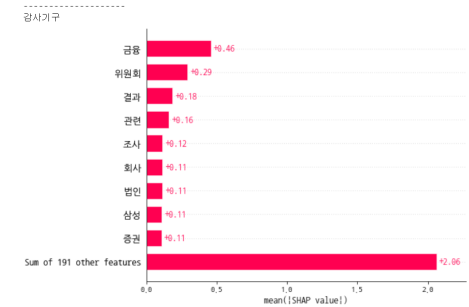
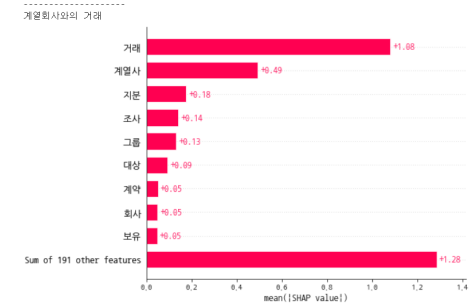
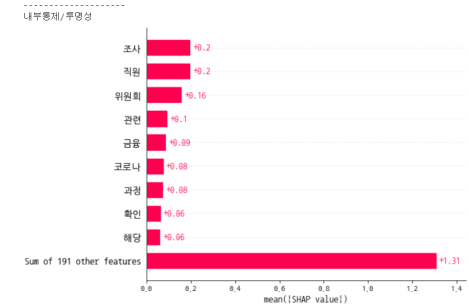
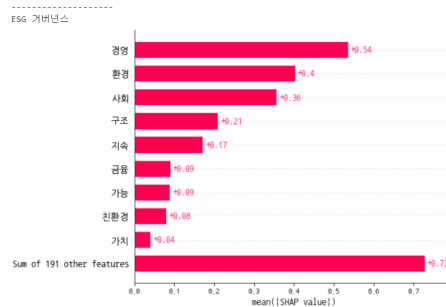
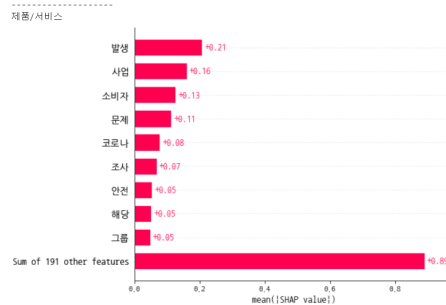
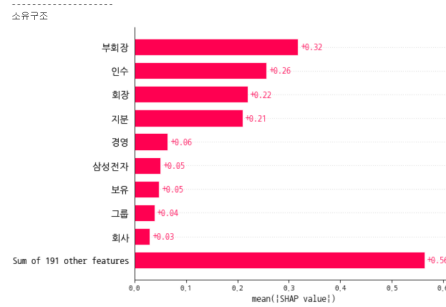
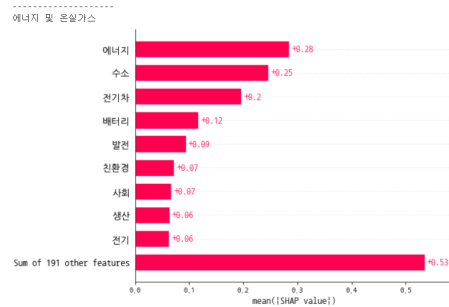
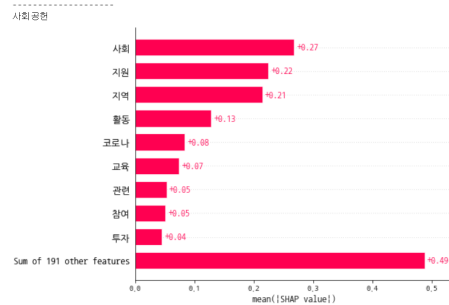
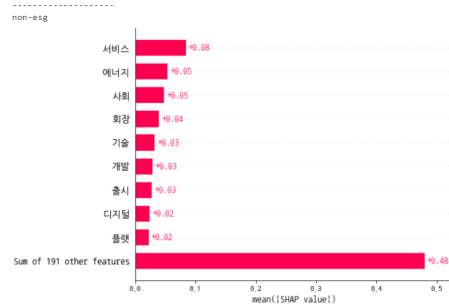
Minor classification performance



4. 모델 성능 및 산출물 : 카테고리 분류

모델 Interpretability & Explainability test : Class 별 단어 중요도

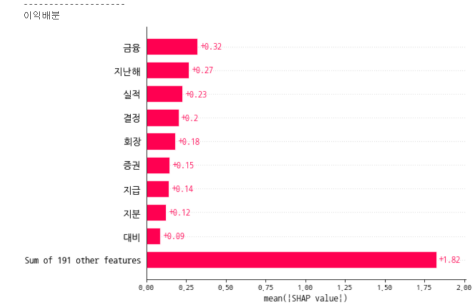
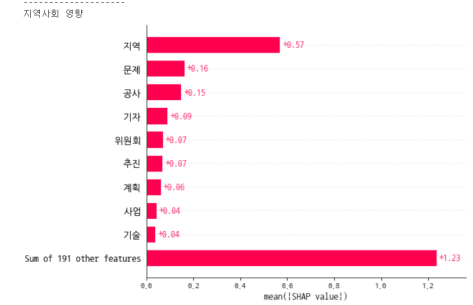
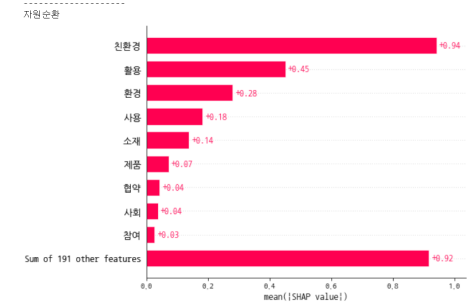
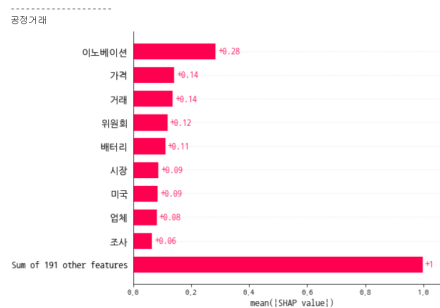
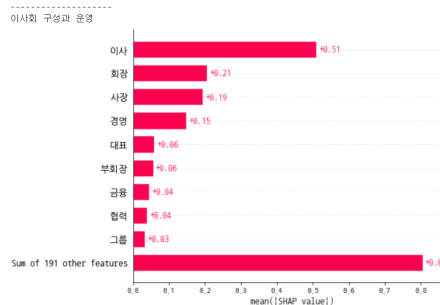
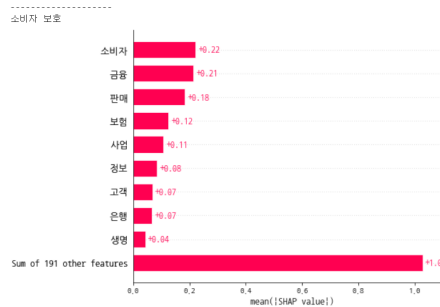
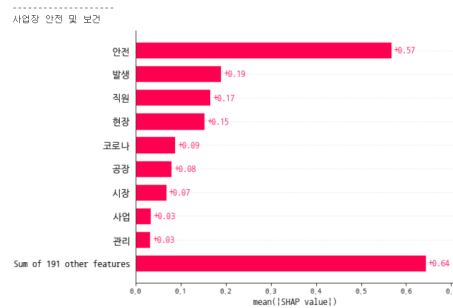
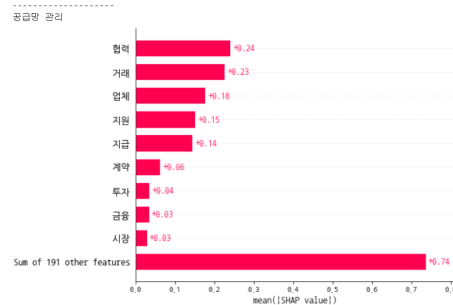
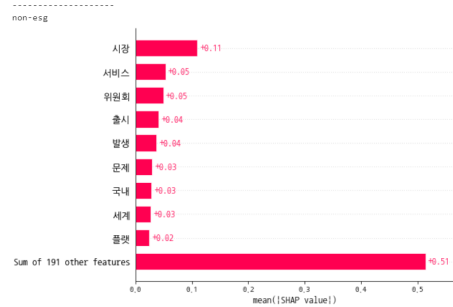
Major classification



4. 모델 성능 및 산출물 : 카테고리 분류

모델 Interpretability & Explainability test : Class 별 단어 중요도

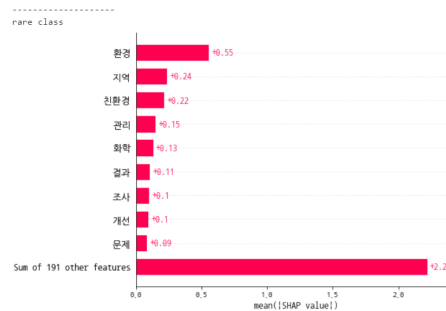
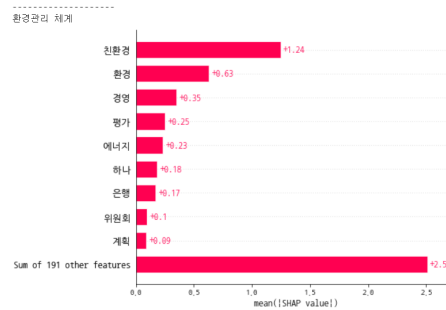
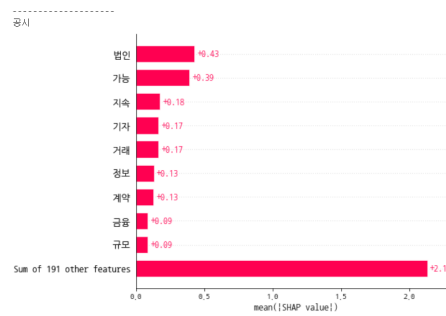
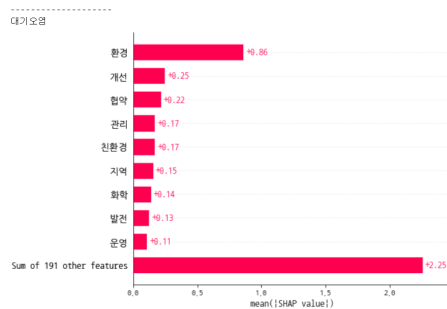
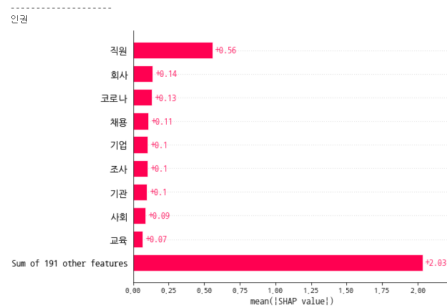
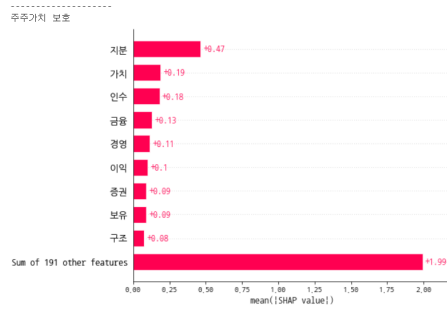
Minor classification



4. 모델 성능 및 산출물 : 카테고리 분류

모델 Interpretability & Explainability test : Class 별 단어 중요도

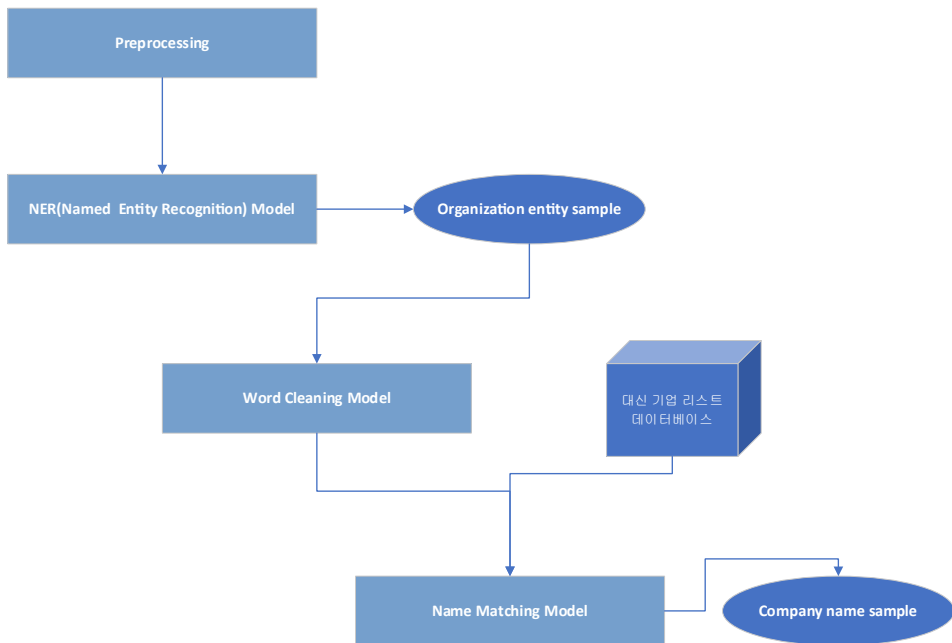
Minor classification



4. 모델 성능 및 산출물 : 기업명 추출

모델 구성도 및 예시

모델 구성도



Hit ratio : 73%

결과 예시

추출된 기업명 (요약본)	추출된 기업명	제목
CJ그룹, CJ/CGV, CJ푸드빌	CJ푸드빌	CJ푸드빌 1년 만에 수장 교체...정성필철 '적자 개선' 발등에 불
KB금융그룹, 대한카누연맹	대한민국카누국가대표팀, KB금융	KB금융, 대한민국 카누 국가대표팀 공식 후원
KT&G, 필립모리스	필립모리스	필립모리스, 필연형 전자담배 시장 1위 유지할까 - 서울경제
LG상, LG이노텍(주), LG	LG	40세 구광모, 재계 4위 LG 회장 등극...구본준 "소임 끝났다 ...
LG화학, 삼성SDI	LG화학	삼성SDI, LG화학, 수소차에 맞서 전기차배터리 경쟁력 확보 서둘러
SK이노베이션	SK	차세대 NCM811 배터리 SK가 먼저 상용화...세계 최초 양산 스타트

예시 엑셀 파일

	A	B	C	D	E	F	G
		id	date	keyword	ner_t	title	
1	0	160053	2018-07-01	CJ대한통운	CJ푸드빌	CJ푸드빌 1년 만에 수장 교체...정성필철 '적자 개선' 발등에 불	
2	1	160276	2018-07-01	KCC	대한민국카누국가대표팀, KB금융	KB금융, 대한민국 카누 국가대표팀 공식 후원	
3	2	160338	2018-07-01	LF	필립모리스	필립모리스, 필연형 전자담배 시장 1위 유지할까 - 서울경제	
4	3	160409	2018-07-01	LG상철건설	LG	40세 구광모, 재계 4위 LG 회장 등극...구본준 "소임 끝났다 ...	
5	4	160414	2018-07-01	LG상철건설	LG화학	삼성SDI, LG화학, 수소차에 맞서 전기차배터리 경쟁력 확보 서둘러	
6	5	160820	2018-07-01	SK케미칼	SK	차세대 NCM811 배터리 SK가 먼저 상용화...세계 최초 양산 스타트	
7	6	161509	2018-07-01	삼성바이오로직스	삼성현대, 한	무늬만 '공익법인'...삼성현대자·한진 등 중수지배 진위부터	
8	7	161535	2018-07-01	삼성바이오로직스	삼성현대, 중앙일보	금융규제, 삼성현대자 시름에 빠뜨린다...중앙일보	
9	8	162144	2018-07-01	구글	삼성현대자, 중앙일보	구글이 '구글'이 '조용히' 삼성현대자...중앙일보	
10	9	162201	2018-07-01	포스코인하철강	국세청	한승희 1년 국세청 개학...국세청 2년차 등락 확정	
11	10	162204	2018-07-01	포스코인하철강	LG	40세 구광모 LG 회장 등 보좌관 부회장 6인방은?	
12	11	162211	2018-07-01	통산	LG유플러스, 에어서울	에어서울, LG유플러스 해로로 시 '일거양득'	
13	12	162531	2018-07-01	한화	한화	한화 최진욱 시장 승진...커뮤니케이션위원장	
14	13	160018	2018-07-02	CJ CGV	CJ ENM	CJ ENM 출범, 미디어서비스 경쟁 시작된다	
15	14	160097	2018-07-02	CJ대한통운	CJ푸드빌	정성필 CJ푸드빌 대표, 애플인자 '해리시' 앞길	
16	15	160124	2018-07-02	DB하이텍	CJ푸드빌	장대영 CJ푸드빌 대표, 애플인자 '해리시' 앞길	
17	16	160146	2018-07-02	DB하이텍	보통연구원	보통연구원 '국내 보형물, 보형물' '보형물' '보형물' '보형물' ...	
18	17	160155	2018-07-02	GKL	보통연구원	태풍 피해복구 예산, 순이익을 늘릴까? 순이익을 늘릴까? ...	
19	18	160185	2018-07-02	GKL	보통연구원	건강 나이 따지는 보형물상 나온다	
20	19	160203	2018-07-02	GS	리치엔코	리치엔코, 보형물상 '국지' 업계 최초 100만 다운로드 돌파	
21	20	160248	2018-07-02	JW중외제약	HDC현대산업개발	HDC현대산업개발, 42억 규모 자부담 자금	
22	21	160252	2018-07-02	JW중외제약	HDC현대산업개발	HDC현대산업개발, 한국농어촌공사·HDC현대산업개발, 한국 농어촌을 위한 업무협약 체결	
23	22	160267	2018-07-02	JW중외제약	HDC현대산업개발	HDC현대산업개발, 한국농어촌공사·HDC현대산업개발, 한국 농어촌을 위한 업무협약 체결	
24	23	160319	2018-07-02	LF	KT&G	KT&G 신세계 순이익, 보형물상 개발 나선다	
25	24	160435	2018-07-02	LG유플러스	LG	구광모 LG 회장, 조영남 취임...'개성' 변화 '상조'	
26	25	160461	2018-07-02	LG유플러스	LG디스플레이, LG	구광모 LG 회장 첫 시범하는 LG디스플레이 '정성화'	