

# kohate 통합 분석 리포트

데이터 드리프트 분석 보고서

생성일시: 2025년 08월 19일 10시 12분

## Dataset Information & Statistics

### Test Dataset

#### Preview

comments	contain_gender_bias	bias	hate
송중기 시대극은 믿고본다. 첫회 신선하고 좋았다.	False	none	none
지현우 나쁜놈	False	none	offensive
알바쓰고 많이 만들면 되지 돈 욕심 없으면 골목식당 왜나 온겨 기댕기게나하고 산에 가서 팔아라	False	none	hate
설마 ♂ 현정 작가 아니지??	True	gender	hate
이미자씨 송혜교씨 돈이 그리 많으면 탈세 말고 그돈으로 평소에 불우이웃에게 기부도 좀 하고 사시죠.	False	none	offensive
일베충들 냐드거리는 것 봐라 짜짜	False	others	hate
아이즈원 힘내세요... 일본 진출도 했으니 일본서 좋은 모습 보여줘도 팬들은 응원 합니다.	False	none	none
강부자 선생님 전미선 비보에 오열을 하셨다니 눈물이 나네요 힘내세요	False	none	none
알았어 그만	False	none	offensive
이영자씨는 진정성 있는 거라면 녹화 불참으로 끝내지 말고 자진하차해라 시청자는 고려도 안하고 일방적 불참은 아닌듯 엠비씨도 시청률 좋아서 고민하는 거 같은데 결방할 게 아니고 폐지해라	False	none	offensive

#### Description

comments	contain_gender_bias	bias	hate
471	471	471	471
471	2	3	3
송중기 시대극은 믿고본다. 첫회 신선하고 좋았다.	False	none	offensive
1	404	342	189

## Info

Non-Null Count	Null Count	Dtype
471	0	object
471	0	category
471	0	category
471	0	category

## Train Dataset

### Preview

comments	contain_gender_bias	bias	hate
○정보부에서 조사하나봐요 .. 승리 곧 자살하겠네 아니..경찰이 질소넣고 자살했다고 하겠네..	False	others	hate
어차피 속바지인데 뭔 노출ㅋㅋ	True	gender	hate
손담비 나훈에 고정으로 가자~~매력 넘침~	False	none	none
개들보네 ㅋㅋㅋㅋㅋ남자 뭔 프로게이머같이생겼네	False	none	hate
누가봐도 한거다. ㅎㅎㅎ	False	none	offensive
처음 본 여자한테 찝적 거렸는데, 여자가 기분 나쁘다 ? 이게 죄야 ? 성적 발언을 한 것도 아니고, 상대가 기분이 나쁘면 무조건 피해자로 규정하나 ? 이건 아니지 싶다 ㅋㅋㅋ	True	gender	none
에이핑큰줄알았더니 짹퉁김남주였네	False	none	offensive
개쳐늙으면 꼴보기 싫어지는게 사람의 마음이지 얘네 야 광고수익으로 번게 많으니 얼굴에 돈들이고 관리하니 금슬이 당연히 좋을수밖에 사람이란 다 그런거임 아무리 이쁘고 잘생겼어도 관리안하면 상대쪽에서 싫어하게되어있음	False	others	hate
미투 최대 피해자김생민 ..불쌍해	True	gender	none
솔직히 인성이 종아리에서 나왔음	False	none	offensive

### Description

comments	contain_gender_bias	bias	hate
6316	6316	6316	6316
6316	2	3	3
○정보부에서 조사하나봐요 .. 승리 곧 자살하겠네 아니..경찰이 질소넣고 자살했다고 하겠네..	False	none	none

## Info

Non-Null Count	Null Count	Dtype
6316	0	object
6316	0	category
6316	0	category
6316	0	category

## Validation Dataset

### Preview

comments	contain_gender_bias	bias	hate
제발~~부끄러운줄알았음..이혼은왜하셨나요??사생 활지저분한쇼호스트가 티비나와서할소린아닌줄..	False	others	hate
윌리엄 그동안 스킁했었는데 오늘 이모습하고 서준이에서 눈여겨볼 아이로 성장했음	False	none	none
우익 꺼지셈	False	others	offensive
이뻐보이진 않는데..	False	none	none
와ㅜㅜ 고배우 연기 진짜ㅜㅜ 와 소름돋음	False	none	none
방금 유동근 아래 회장여자 한테 엄포 놓을때 포스 ㅎㄷㄷ간만에 정도전때 포스 나오더만 진작 이렇게 만들지 이상하게 캐릭터를 만들어 놓아서.	False	none	none
백상까지 조작한 sm.문재인들어서고 방송 문화 경제 등 비리와 부패가 더 심각하고 망하고있다.조작사기다 sm비리와 조작으로 인기구걸하는데도 조사없는 문재인끌어내려야지	False	others	hate
당신을 보면 딱 이구절만 생각나네요바보가 신념을 가지면 무섭습니다.결혼한다니 축하드리고 잘사시길 바랍니다.자기 주제에 맞게 사시길.	False	none	offensive
육아하는 아빠라니 육아는 부모라면 당연한거고요....ㅋ..;;	False	none	none
저런곳에 일하는여자가 저런거 당했다고 하루하루 괴롭게 지낸다고?개가 웃겠다ㅋㅋ돈만주면 2차 나가는 거 좀 좋아하는 애들 아니냐? 그리고 성매매 업소도 아니라고? 성매매 업소도 성매매 업소라고 안하고 성매매하거든?어이없네	True	gender	hate

### Description

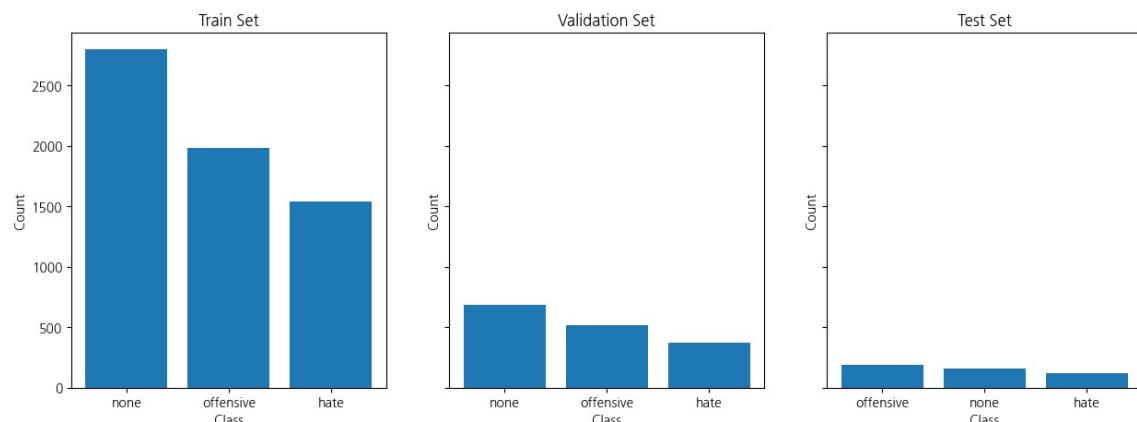
comments	contain_gender_bias	bias	hate
1580	1580	1580	1580
1580	2	3	3
제발~~부끄러운줄알았음..이혼은왜하셨나요??사생활지저분한쇼호스트가 티비나와서할소린아닌줄..	False	none	none
1	1328	1022	688

## Info

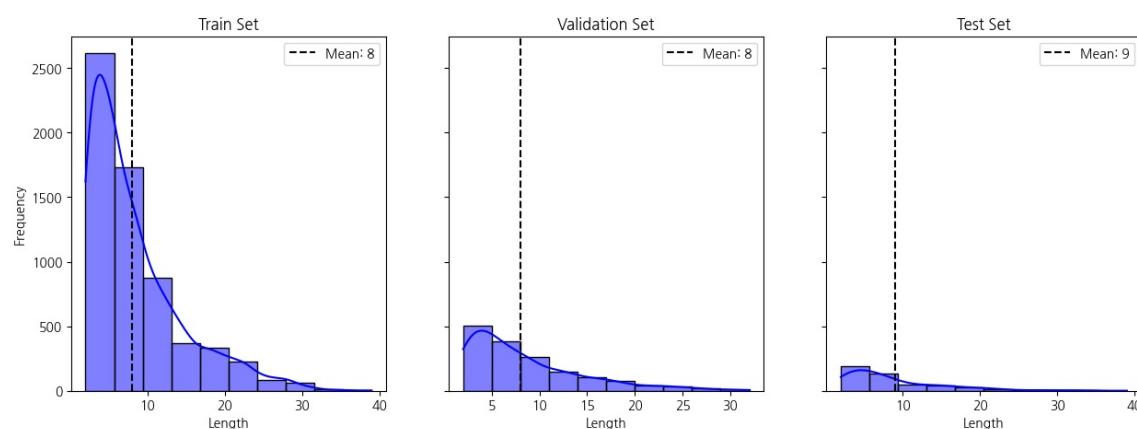
Non-Null Count	Null Count	Dtype
1580	0	object
1580	0	category
1580	0	category
1580	0	category

## Visualizations

### Class Distribution



### Document Length Distribution



Dataset	Longest Sentence	Shortest Sentence	Mean Sentence Length	Sum of Sentences
Train	39	2	8	6316
Validation	32	2	8	1580
Test	39	2	9	471

## Word Cloud

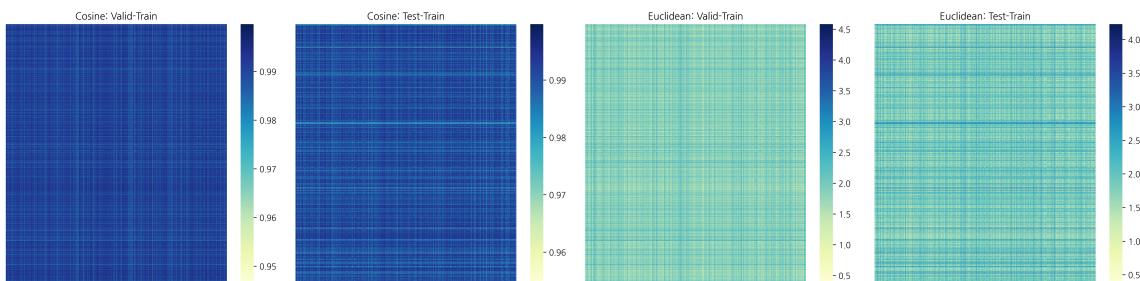


## Data Drift Analysis Results

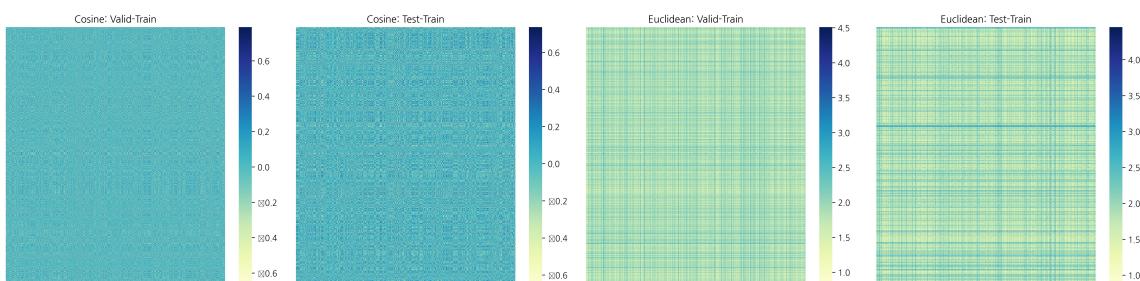
Train: (6316, 768), Valid: (1580, 768), Test: (471, 768)

**PCA Reduced Dimension:** 400.0

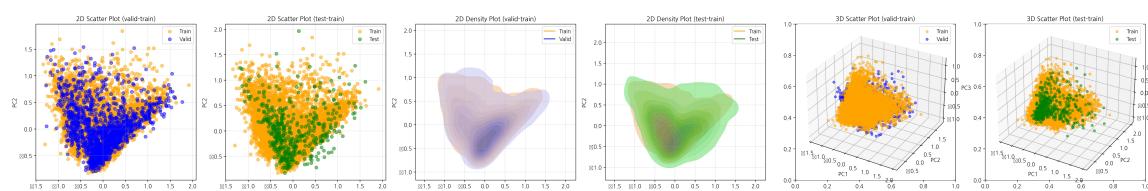
## Embedding Distance (Original Dimension)



## Embedding Distance after PCA



## Embedding Visualization after PCA



## Drift Score Summary

- MMD: score = 0.0000, drift = False
- Wasserstein Distance: score = 0.5850, drift = True
- KL Divergence: score = 0.4175, drift = True
- JensenShannon Divergence: score = 0.5225, drift = True
- Energy Distance: score = 0.0025, drift = False

## LLM Explanation

## 드리프트 분석 해설 \*\*[기술적 분석]\*\* 분석 결과, Wasserstein Distance (0.5850), KL Divergence (0.4175), JensenShannon Divergence (0.5225)의 값이 모두 0에 가까운 값보다 상당히 높게 나타났습니다. 이는 데이터 분포가 크게 변화하고 있다는 것을 의미하며, 모델이 현재 학습하는 데이터와 과거 학습 데이터 간의 차이가 상당하다는 것을 시사합니다. 특히, Wasserstein Distance는 데이터 분포의 유사성을 측정하는 데 사용되므로, 0.5850이라는 높은 값은 모델의 성능에 심각한 영향을 미칠 수 있음을 나타냅니다. \*\*[현 상황 분석]\*\* 현재 드리프트 상황은 모델의 학습 성능에 심각한 위협이 됩니다. 높은 드리프트 지표는 모델이 새로운 데이터에 적응하지 못하고, 예측 정확도가 떨어지거나 심지어는 성능이 급격히 저하될 수 있습니다. 특히, 모델이 텍스트 데이터에 특화되어 있다면, 데이터의 내용이나 스타일이 변화함에 따라 모델의 성능이 크게 저하될 가능성이 높습니다. \*\*[시각적 분석]\*\* PCA 시각화 결과가 제시되지 않았으므로, 현재 데이터 분포 변화를 정확하게 시각적으로 해석하는 것은 불가능합니다. 그러나 제시된 드리프트 지표들을 고려했을 때, 데이터의 특징 (feature) 공간에서 데이터 분포가 크게 변하고 있음을 추정할 수 있습니다. PCA 시각화를 통해 데이터 분포 변화의 패턴을 확인하고, 어떤 특징(feature)이 변화의 원인인지 파악하는 것이 중요합니다. \*\*[권장사항]\*\* 즉시 모델 재학습을 고려해야 합니다. 현재 드리프트 지표가 높게 유지된다면, 모델의 성능이 지속적으로 저하될 가능성이 높습니다. 또한, PCA 시각화를 통해 데이터 분포 변화의 패턴을 분석하고, 변화의 원인이 되는 특징 (feature)를 파악하여 모델을 재학습할 때 해당 특징에 대한 가중치를 조정하는 것이 좋습니다. 지속적으로 드리프트 지표를 모니터링하고, 변화가 감지되면 즉시 대응하는 시스템을 구축하는 것이 중요합니다.

DataDrift Dataclinic System

@KETI Korea Electronics Technology Institute, 2025