

LLM 평가 프레임워크 비교 분석

→ 각 프레임워크의 특징이 무엇이고, 어떤 장단점 평가를 수행하는 지 파악한 다음에 사용할 평가 프레임워크를 선택하고 사용해야 한다.

[비교해볼 대표 LLM 평가 프레임워크]

- 1) RAGAS, DeepEval 등 'LLM-as-a-judge' 기반 프레임워크
- 2) LangSmith, TruLens 등 observability 프레임워크
- 3) Huggingface evaluate
- 4) ARES
- 5) AutoRAG

※ LLM-as-a-judge 프레임워크

→ "LLM이 평가자가 되어 평가한다"라는 concept

(원리) → 대변 "프롬프트 엔지니어링"이다.

1) RAGAS

- 가장 유명한 평가 프레임워크
- Retrieve 한 단락이 질문과 얼마나 관련 있는지를 LLM으로 평가
 - 단락 내용을 모두 LLM에게 전달해야 하기 때문에 '굉장히 비싸다'
 - RAGAS는 'RAG 평가'에 사용하기는 비싼 편, LLM 평가에는 많이 사용됨.
- LLM 평가 뿐 아니라 여러 평가 Metrics로 측정하는 장점을 가짐
 - Agent의 Tool 사용에 대한 평가도 포함 (RAGAS에만 거의 원천코딩 포함)
 - SQL Metric과 정통 NLP 스코어보드 포함
 - Vector DB가 아닌 RDB 기반 DB를 LLM Application에 사용시 LLM이 SQL Query 문을 작성하도록 하는 경우도 존재하는데, 이 때, 이를 평가하는 역할.
- 최근 자연 X (영어간 자연), 한자간 성능이 높아질 수 있음.

2) DeepEval

- RAGAS, G-Eval 등 LLM-as-a-judge 메트릭 포함
- AI Safety 관련 메트릭 포함 (LLM-as-a-judge, RAGAS는 X)
- 여러 벤치마크 데이터셋 지원, CI/CD 통합 지원
 - ↳ LLM 평가 만 필요하거나,
 - ↳ 지속적으로 평가를 자동으로 수행하는 각종 서면
- ↳ 관련 LLM application이 기존 데이터셋은 활용되도 충분히 General한 도메인 경우
- 다국어 지원 X
- Retrieval Metric X → RAG를 전혀 고려하고 있지 X

3) OpenAI Evals

- OpenAI에서 직접 제작한 평가 프레임워크
- OpenAI Dashboard에서 사용 가능
- LLM 평가에 집중되어 있음.
- OpenAI 도입한 사용 가능.
 - 즉, OpenAI 도입이 아닌 도입을 도입한 경우 사용이 제한됨.

4) LangSmith

- LangChain과 연동되는 모니터링 및 디스토크 플랫폼.
- 기본적으로는 모니터링 등, 평가 기능 지원
- 'LLM-as-a-judge'를 기본적으로 제공, 커스텀 메트릭 지원
- 평가 결과를 쉽게 확인 가능한 대시보드 지원
- 치명적 반영 : 부분적 유료...

5) TruLens ; LangSmith와 비슷하지만 조금 더 평가에 초점을 맞춘 프레임워크

- Human-in-the-loop 등 반복적 평가에 적합
 - 평가 과정에 사람이 개입할 수 있도록 해 신뢰성↑
- 결과 확인을 위한 대시보드 지원
- AI Safety 관련 메트릭 포함
- Observability 등의 성격이 있는 평가 프레임워크

6) Hugging-face Evaluate

- 프레임워크, 라이브러리 0 → 평가 메트릭을 모아둔 것.
→ 평가를 ML 메트릭과 데이터셋을 포함한 라이브러리
- 정통 NLP 메트릭 다수를 쉽게 사용할 수 있도록 제공
→ 'LLM-as-a-judge'보다는 전통적인 NLP 평가 메트릭을 많이 써임
- 간단한 사용법으로 빠르게 메트릭 사용 가능
- 다른 평가를 라이브러리들에 의존 관계가 없음
→ 다른 라이브러리를 Worrying 하찮은 것이 아님

7) ARES

- 스탠포드 Univ.에서 만든 LLM평가 프레임워크
- 평가를 위한 SLLM을 파인튜닝 하는 것이 주임
→ 가장 높은 성능 비용이 필요하지만, 'LLM-as-a-judge'중 정확도가 높은 편
→ 하지만 파인튜닝이 필요하다는 장이 가장 큰 단점.
→ Local 모델로 사용해야하기 때문에 GPU도 필요.

8) AutoRAG

- RAG에 특화된 최적화 프레임워크
RAG 최적화에 초점
- NLP 메트릭, LLM-as-a-judge 등 여러 메트릭 포함
→ ARES를 제외한 모든 프레임워크가 통합되어 있음.
- Retrieval 평가 및 최적화 서원 (Retrieval GT 기반의 정확한 Retrieval 평가 자명)
- 평가 데이터셋 생성에서 한국, 영어, 일본어 서원
- 평가 마시보드 서원

※ 평가 프레임워크 정리

최적화

Dashboard

Observability

→ 기능.

LLM-as-a-judge

선형적 NLP Metrics

기타 평가 메트릭

→ 메트릭.