

Retrieval 평가와 메트릭

Retrieval Metric

- ① Rank-Unaware (순서를 고려하지 않는 Metric)
→ Recall, Precision, F1
- ② Rank-aware (순서를 고려하는 Metric) : 전통적인 검색 시스템에서 사용되는 방법
→ mAP, mRR, NDCG

평가 방식 → 'Retrieval contents'와 'Retrieval GT'의 ID를 비교해서 측정하는 방식 채택

↓

검색 결과

↓

정답

→ 즉 정답에 존재하는 ID를 가진 내용 같 지만 있는 것 같아 보이는 것.

		실제 정답	
		True	False
분류 결과	True	True Positive	False Positive
	False	False Negative	True Negative

○ Precision

→ Retrieved contents 중 정답 수
Retrieved contents의 전체 수 → 검색해온 전체 수 (= TopK)

→ 'Positive 정답률', 'PPV (Positive Predictive Value)'라고도 불림

즉, True라고 여쭙한 것 중, 실제 True인 것의 비율

○ Recall (Sensitivity, hit rate)

→ GT 중 맞은 정답 수 → 검색 결과에 포함된 GT 수
Retrieved GT의 전체 수 → 전체 GT 수

→ 실제 True인 것 중, True라고 여쭙한 것의 비율

→ 'Auto RAG' 프레임워크에서는 GT가 기본적으로 2바이트 255로

ex) Retrieved Contents (ID)

test-1 0
Prod-1 x
test-2 0
Prod- x

Retrieval GT

[test-1, test-2]
0 0
[test-3]
x

→ Retrieval 평가에서는 '정답을 가져온 것'이 '가
바뀌 필요하기 때문에 Recall도 많이 봄

→ 정답을 가져오면 점수를 한 수가 x 때문

→ Top-k를 ↑ 하면 Recall도 증가하지만

그만큼 다른 공짜도 같이 가져와서 '정답' 확률이 ↑

「여기 ID가 1바이트 255
이하 ID는 정답이」

test-1 or test-2
AND
test-3

$$\Rightarrow \frac{\text{GT 중 정답 개수}}{\text{전체 GT 수}} = \frac{1}{2} = 0.5$$

0 F1 Score

→ 'Precision'과 'Recall'은 각각은 반반씩
이들 보정하여 커서 F1-Score '상승'

$$\rightarrow 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (\text{Precision과 Recall의 조화-평균})$$

→ 0이 사이의 값. 1에 가까울 수록 성능이 좋다.

※ Top-k로 보고 싶을 때, 정답과 일대일 오답을 바꿔서

→ 1, 2, ... 보라 50까지 Top-k를 늘리면 Recall이 증가함 (거의 1까지 도달)

→ 즉, 많이 가져다서 정답이 포함될 수 있도록 한 다음,
결과를 ReRanker에 넣어서 정답 결과로 조정.

⇒ 이렇게 하얏해 보자.

0 MRR (Mean Reciprocal Rank) ; Mean RR

→ 정답 단락 (Passage) 가 몇 번째로 나왔는지 (RR)에 대한 평균

ex) Retrieved Contents (ID)

ID-1 X
ID-2 0
ID-3 X
ID-4 0

Retrieved GT

ID-2
ID-4

① (GT ID-2의 RR)

→ Contents에서 2번째 순서 등장

$$\rightarrow \frac{1}{\text{반정된 수의 개수의 수}} = \frac{1}{2} = 0.5$$

② (GT ID-4의 RR)

→ Contents에서 4번째 순서 등장

$$\rightarrow \frac{1}{4} = 0.25$$

③ RR의 평균 (MRR) 계산

$$\frac{0.5 + 0.25}{2} = 0.375$$

o mAP (Mean Average Precision) ; Rank-Unaware의 Precision과 다른 Precision

→ AP (Average Precision)의 평균 (평균의 평균)

step 1) 각 gt 당 precision 계산 (Precision)

step 2) 위들의 precision의 평균 계산 (Average Precision)

step 3) AP의 평균 계산 (mAP)

ex)

Retrieved Contents (ID)	
test-1	o
Pred-1	x
test-2	o
Pred-	x

Retrieval GT	
[test-1, test-2]	o
[test-3]	x

step 1)

[test-1, test-2] → ① GT 'test-1'이 첫번째로 등장
⇒ $\frac{1}{1} = 1$

② test-1이 첫번째, test-2가 세번째 등장

step 2) $\frac{1 + 2/3}{2} = \frac{5}{6}$

$\frac{0}{4} = 0$

step 3) $\frac{5/6 + 0}{2} = \frac{5}{12}$ (mAP)

⇒ $\frac{\text{포함된 GT 수}}{\text{반정되기까지 수}} = \frac{2}{3} = 0.66$

o NDCG (Normalized Discounted Cumulative Gain)

→ 상위 순위에 더 높은 가중치를 부여, 결과의 순위 품질을 0으로 평가

- step 1) CG (Cumulated Gain) → 관련성의 누적 합, 순위고려 X ($CG_p = \sum_{i=1}^p rel_i$)
Andorid에서는 0/1 만 사용 (관련성 점수 제외)
- step 2) DCG (Discounted Cumulated Gain) → 관련성 점수를 소위 나열한 값 ($DCG_p = \sum_{i=1}^p \frac{rel_i}{\log_2(i+1)}$)
값이 작아지면 증가, 관련성 영향력이 후순위로 갈수록 ↓
- step 3) IDCG (Ideal Discounted Cumulated Gain) → DCG의 가장 이상적인 값 (순서 정렬 결과와는 상관)
 → 즉 검색 순위 결과가 가장 좋은 경우 DCG
- step 4) NDCG (Normalized Discounted Cumulative Gain) → $\frac{DCG}{IDCG}$, 이상적 상태에 얼마나 가까운가?
 0이서이, 1에 가까운수록 좋음.