

ZTM RAG: RAG Fundamentos

terça-feira, 1 de julho de 2025 13:01

Sobre LLMs (Large Language Models)

Key Methods:

- Zero Shot Learning: Modelo prevê a resposta apenas com uma descrição da pergunta.
- One-shot Learning: Além da descrição do problema, é dado um exemplo do que o modelo deve fazer.
- Few-shot Learning: Modelo recebe a descrição do problema e alguns exemplos do que deve fazer.

Em fine tuning, geralmente é usado o Few-shot Learning. Isso geralmente é bem caro e não usado mais.

Por isso, fine tuning se tornou legacy e agora os modelos usam centenas de bilhões de parâmetros para saber como resolver os prompts sem precisar de fine tuning.

Essa quantidade de dados gerou um viés de sentimento em tarefas raciais por não terem limpado os dados corretamente.

Aumentar o tamanho do modelo, o tamanho do dataset e o poder computacional melhora o resultado, mas também aumenta o custo.

As LLMs conseguiram ser melhores que os humanos quando implementaram o método SAFE (Self Attention Fine Tuning and Evaluation), que permite os modelos verificarem a acurácia do próprio conteúdo.

RAG

O que é RAG?

É um modelo híbrido que melhora a geração de texto usando informações de documentos, levando a respostas mais precisas.

Componentes principais de uma arquitetura básica de RAG

- **Input Query**
A pergunta do usuário ou prompt que inicia os outros processos.
- **Retriever**
Pega documentos ou passagens relevantes de um grande corpus baseado no input query. Isso gera os **Retrieved Documents**, que proveem o contexto e informação necessária para gerar uma resposta precisa e relevante.
- **Generator**
Cria uma resposta baseada na input query e nos retrieved documents, usando LLMs.

FAISS

É uma biblioteca criada pelo Facebook para busca eficiente de similaridade entre vetores de alta-dimensão. Como um bibliotecário muito bom numa biblioteca gigante.

É bom pela velocidade, escalabilidade e flexibilidade.

Como funciona: Indexando dados em estruturas otimizadas para busca rápida de similaridade.

Tipos de FAISS Indexes

- **Flat Index**
Simple método de força bruta que checa todos os vetores. Ideal para pequenos datasets.
- **IVF** (Inverted File)
Divide os dados em clusters para busca rápida. Ideal para grandes datasets.
- **HNSW** (Hierarchical Navigable Small Worlds)
Usa um graph base aproch para buscas bem rápidas. Ideal para aplicações de tempo real.

Generation Model Parameters

São settings que controla o comportamento do modelo de geração de texto. Modificar esses parâmetros podem melhorar significamente a performance e a qualidade dos modelos. Entre esses parâmetros, estão temperature, top-k sampling, top-p sampling, repetition penalty, sampling mode.

- **Temperature**
Controla a aleatoriedade do texto gerado, balanceando entre outputs focados e criativos. Valores menores são mais focados e valores maiores são mais aleatórios e criativos. Valor default é "1.0".
- **Top-K Sampling**
Limita a escolha da próxima palavra possível para as mais prováveis.
Ex: Top-K = 50 = apenas as 50 palavras mais prováveis serão consideradas na hora de gerar a próxima palavra.
- **Top-p (Nucleus) Sampling**
Limita a gama de palavras consideradas para a próxima palavra a gerar.
Se "top_p = 1.0", significa que 100% das palavras são consideradas para a próxima palavra.
É semelhante ao Top-K.
- **Repetition Penalty**
Reduz as frases repetitivas, fazendo respostas mais diversas e humanas.
Valor padrão = 1.0 - Acima disso ele penaliza e abaixo ele recompensa repetição.
- **Sampling Mode**
Adiciona mais aleatoriedade, criando textos mais variáveis e criativos.

LongRAG

LongRAG utiliza mais tokens do que a forma tradicional de rag para gerar contexto para as respostas, reduzindo a fragmentação. RAG normalmente utiliza ~100 tokens de contexto. LongRAG utiliza ~4000 tokens (que da aproximadamente 3 mil palavras), fornecendo um contexto muito maior, gerando uma resposta muito mais completa.
É fácil de aplicar.

LightRAG

LightRAG combina RAG tradicional com índice de grafo de conhecimento, integrando tanto componentes textuais quanto estruturais (entidades e relações). Ideal para dados complexos. Não é tão simples de aplicar.

LightRAG vs Outros RAGs em diferentes datasets

Table 1: Win rates (%) of baselines v.s. LightRAG across four datasets and four evaluation dimensions.

	Agriculture		CS		Legal		Mix	
	NaiveRAG	LightRAG	NaiveRAG	LightRAG	NaiveRAG	LightRAG	NaiveRAG	LightRAG
Comprehensiveness	32.69%	<u>67.31%</u>	35.44%	<u>64.56%</u>	19.05%	<u>80.95%</u>	36.36%	<u>63.64%</u>
Diversity	24.09%	<u>75.91%</u>	35.24%	<u>64.76%</u>	10.98%	<u>89.02%</u>	30.76%	<u>69.24%</u>
Empowerment	31.35%	<u>68.65%</u>	35.48%	<u>64.52%</u>	17.59%	<u>82.41%</u>	40.95%	<u>59.05%</u>
Overall	33.30%	<u>66.70%</u>	34.76%	<u>65.24%</u>	17.46%	<u>82.54%</u>	37.59%	<u>62.40%</u>
	RQ-RAG	LightRAG	RQ-RAG	LightRAG	RQ-RAG	LightRAG	RQ-RAG	LightRAG
Comprehensiveness	32.05%	<u>67.95%</u>	39.30%	<u>60.70%</u>	18.57%	<u>81.43%</u>	38.89%	<u>61.11%</u>
Diversity	29.44%	<u>70.56%</u>	38.71%	<u>61.29%</u>	15.14%	<u>84.86%</u>	28.50%	<u>71.50%</u>
Empowerment	32.51%	<u>67.49%</u>	37.52%	<u>62.48%</u>	17.80%	<u>82.20%</u>	43.96%	<u>56.04%</u>
Overall	33.29%	<u>66.71%</u>	39.03%	<u>60.97%</u>	17.80%	<u>82.20%</u>	39.61%	<u>60.39%</u>
	HyDE	LightRAG	HyDE	LightRAG	HyDE	LightRAG	HyDE	LightRAG
Comprehensiveness	24.39%	<u>75.61%</u>	36.49%	<u>63.51%</u>	27.68%	<u>72.32%</u>	42.17%	<u>57.83%</u>
Diversity	24.96%	<u>75.34%</u>	37.41%	<u>62.59%</u>	18.79%	<u>81.21%</u>	30.88%	<u>69.12%</u>
Empowerment	24.89%	<u>75.11%</u>	34.99%	<u>65.01%</u>	26.99%	<u>73.01%</u>	45.61%	<u>54.39%</u>
Overall	23.17%	<u>76.83%</u>	35.67%	<u>64.33%</u>	27.68%	<u>72.32%</u>	42.72%	<u>57.28%</u>
	GraphRAG	LightRAG	GraphRAG	LightRAG	GraphRAG	LightRAG	GraphRAG	LightRAG
Comprehensiveness	45.56%	<u>54.44%</u>	45.98%	<u>54.02%</u>	47.13%	<u>52.87%</u>	51.86%	<u>48.14%</u>
Diversity	19.65%	<u>80.35%</u>	39.64%	<u>60.36%</u>	25.55%	<u>74.45%</u>	35.87%	<u>64.13%</u>
Empowerment	36.69%	<u>63.31%</u>	45.09%	<u>54.91%</u>	42.81%	<u>57.19%</u>	52.94%	<u>47.06%</u>
Overall	43.62%	<u>56.38%</u>	45.98%	<u>54.02%</u>	45.70%	<u>54.30%</u>	51.86%	<u>48.14%</u>

NaiveRAG = RAG padrão

A porcentagem representa o quanto de vezes o modelo ganhou em relação ao outro modelo. Percebe-se que apenas quando os dados são mistos que o LightRAG perdeu e unicamente para o GraphRAG.