



ارزیابی کارشناس هوش مصنوعی | مانا افزار خاور

توضیحات پروژه:

هدف این ترسک، ساخت یک سیستم پایه (RAG) است که متون ورودی را به Embedding Vectors تبدیل کرده، در یک پایگاه داده برداری نگهداری می‌کند و هنگام دریافت پرسش کاربر، با جستجوی شباهت (Similarity Search) و تولید پاسخ توسط LLM، نتیجهٔ نهایی را به همراه منابع ارائه می‌دهد.

مواردی که لازم است انجام شود:

1. انتخاب و توجیه پایگاه داده برداری

- استفاده از یک Vector DB.
- مستندسازی دلیل انتخاب بر اساس: کارایی در جستجوی شباهت و درج داده، مقیاس‌پذیری و شیوهٔ استقرار، پشتیبانی جامعهٔ کاربری و ابزارهای توسعه.

2. مدل‌سازی داده

- طراحی اسکیمای ذخیره‌سازی شامل: Document Text (متن بخش/چانک)، Embedding Vector
- Metadata (شناسهٔ سند، منبع، تاریخ، برچسب موضوعی و ...).
- توضیح چرایی این ساختار و اینکه چگونه به فیلترینگ بهتر، بهبود سرعت بازیابی و ارتقای کیفیت RAG کمک می‌کند.

3. منطق سیستم RAG

◦ *Ingestion Pipeline*

تقسیم متن به چانک‌های مناسب (طول/همپوشانی مشخص)،
تولید Embedding با مدل سبک (مثلًا all-MiniLM-L6-v2)،
ذخیره بردار + متادیتا در پایگاه داده برداری
امکان قرار دادن PDF در یک پوشه مشخص

◦ *Query Pipeline*

دریافت سؤال و تولید Embedding،
بازیابی k چانک مشابه (با امکان فیلتر بر اساس متادیتا)،
ساخت پرامپت حاوی دستورالعمل‌ها + زمینه بازیابی شده،
ارسال به LLM و تولید پاسخ نهایی به همراه منابع.

جداسازی کامل لایه RAG از LLM و استفاده از یک API رایگان یا کم‌هزینه (مثلًا JabirLLM یا Gemini Free) برای تولید پاسخ؛ به‌نحوی که LLM فقط بر پایه context بازیابی شده پاسخ بدهد و از خروج مستقیم محتوای دیتابیس بدون RAG جلوگیری شود.

4. پیاده‌سازی API با FastAPI (بخش اختیاری و انتیازی)

◦ POST /ingest

ورودی: یک سند یا مجموعه‌ای از متون.

عملیات: Chunking، تولید Embedding (مثلًا با sentence-transformers/all-MiniLM-L6-v2) با Chunking یا هر مدل دلخواه) و ذخیره در Vector DB.

خروجی: تأیید ثبت + خلاصه تعداد چانک‌ها/توكن‌ها/ابعاد بردار و شناسه‌ها.

◦ POST /query

ورودی: پرسش کاربر (+ پارامترهای اختیاری مثل k، فیلتر متادیتا).

عملیات: تبدیل پرسش به بردار، Similarity Search، انتخاب k نتیجه‌برتر، ساخت پرامپت ترکیبی و ارسال به LLM.

خروجی: پاسخ نهایی + فهرست context chunks (منابع، امتیاز شباهت).

◦ مستندسازی خودکار API با Swagger/Redoc (FastAPI docs)

5. کیفیت و مستندات

- ثبت تنظیمات کلیدی (ابعاد بردار، k، آستانه شباهت، طول چانک/همپوشانی).
- توضیح طراحی پرامپت برای کاهش hallucination و ادار کردن مدل به استناد به منابع.
- راهنمای اجرا (README): نحوه راهاندازی، متغیرهای محیطی، نمونه درخواست/پاسخ.

نکات کلیدی و موارد مهم

۱- زمان انجام تسک: از زمانی که این تسک دریافت می‌شود، ۷۲ ساعت فرصت دارید تا تسک را انجام داده و ارسال کنید.

۲- محیط اجرا: پروژه را به صورت پابلیک در Github منتشر کنید.

- مخزن باید شامل موارد زیر باشد:

- کد کامل

- فایل environment.yml یا requirements.txt

- مستندات API (در صورت پیاده سازی)

- README.md

۳- کیفیت کد و گزارش:

- کد تمیز، خوانا، مستند و قابل اجرا باشد.

- نتایج و شیوه ارائه به شکل مرتب و شفاف بیان شود.

۴- مصاحبه فنی آنلاین: در جلسه مصاحبه فنی، ممکن است درباره کد و نحوه پیاده سازی شما سؤال شود؛ آماده توضیح منطق و روش کارتان باشید.

۵- ارتباط با ما: در صورت داشتن هرگونه سؤال یا ابهام، می‌توانید از طریق همین ایمیل با ما در ارتباط باشید.

توضیحات اضافی

- می‌توانید برای فیلترینگ متادیتا (مانند تاریخ/منبع/برچسب) در مرحله بازیابی، گزینه‌های اختیاری API قرار دهید.

- توجه به حریم خصوصی: اطمینان از عدم ارسال داده حساس به سرویس‌های بیرونی (در صورت نیاز به سانسور/ناشناس‌سازی).

- تمامی بخش‌های پیاده سازی شده با توابع درج شود.