

قسمت اول

در این پروژه یک موتور جستجو برای یکی از سایت‌های معروف زومیت (<https://www.zoomit.ir/>) یا مجله فرادرس (<https://blog.faradars.org/>) ایجاد می‌کنیم. شما باید یک خزنه بنویسید که از آدرس اصلی داده شده‌ها آغاز به

کار کند و ارجاعاتی که در همین دامنه هستند را دنبال کند و حداقل 100 صفحه را جمع آوری کند.

در سایت زومیت درمنوی فناوری، هر بخش مسیر مشخصی دارد (مثال <https://www.zoomit.ir/computer-learning/>) که اگر صفحاتی از سه مورد از این زیر بخش‌ها را جمع کنید در انجام قسمت دوم تمرینکار شما تسهیل خواهد شد.

در سایت مجله فرادرس هر مقاله یک دسته‌بندی دارد که بالای عنوان به صورت لینک آبی نوشته شده است. در اینجا نیز اگر

مقالتی از سه دسته را جمع کنید در انجام قسمت دوم تمرین کار شما تسهیل خواهد شد.

بخش‌های مختلف مفید هر صفحه مانند عناوین و متن... را جمع آوری کنید. برای پردازش صفحات HTML می‌توانید از کتابخانه‌های مناسب موجود استفاده کنید و محتوای بخش‌های مختلف صفحه را استخراج کنید.

سپس محتوای جمع شده را با استفاده از elastic search یا برنامه خودتان که در تمرین قبل با آن کار کرده‌اید شاخص‌گذاری کنید. با وارد نمودن تعدادی پرس و جوی دلخواه عملکرد مناسب جستجو را نشان دهید

به صورت اختیاری بر اساس بخش‌های مختلف صفحه مکانیزم رتبه‌بندی ایجاد کنید. به عنوان مثال اگر در سندی کلمه در عنوان متن ظاهر شده بود، این سند در رتبه‌بندی نسبت به سندی که کلمه را در متن خود دارد قرار بگیرد. با وارد نمودن تعدادی پرس و جوی دلخواه عملکرد مناسب جستجو را نشان دهید. نحوه پیاده‌سازی خزشگر و موتور جستجو، به همراه انواع پرس و جوها که قابلیت برنامه شما را نشان می‌دهند را در گزارش نهایی ذکر کنید.

قسمت دوم

صفحات مرتبط با حداقل سه بخش مختلف سایت (که در قسمت اول شرح داده شد) را خزش کنید. از هر بخش حداقل 25 صفحه را خزش کنید. سپس برای صفحات هر بخش برچسب همان بخش را در نظر گرفته و یک رده بند بیز ساده را برای

رده‌بندی صفحات اجرا کنید. به این منظور می‌توانید از کتابخانه‌های هر زبانی که استفاده می‌کنید بهره ببرید یا از ابزارهای آماده مثل وکا استفاده کنید. شرح کار و دقت رده بند را گزارش کنید.

به صورت اختیاری می‌توانید صفحات را خوشه‌بندی کنید و خوشه‌های حاصل شده را با برچسب واقعی اسناد مقایسه کرده، معیارهای ارزیابی مناسب را محاسبه کنید.

و در آخر نیاز به یک ویدئو دارم که ضبط بکنید از صفر تا صد پروژه رو توضیح بدید
