

گزارش کار نمایه گذاری

در این پروژه هدف نمایه گذاری تعداد زیادی داده مربوط به کورس های آموزشی میباشد.

نمایه گذاره

برای نمایه گذاری داده ها از اسم پلتفرم آن کورس استفاده میکنیم یعنی هر پلتفرم یک ایندکس در الاستیک سرچ دارد، همچنین برای بهینه سازی سرعت در دسترسی به ارائه دهندگان از نمایه گذاری معکوس به ارائه دهنده استفاده میکنیم زیرا هر کورس میتواند چندین ارائه دهنده داشته باشد یعنی هر پرنتر یک ایندکس از کورس ها در الاستیک سرچ دارد، برای دسته بندی نیز از نمایه گذاره معکوس استفاده می کنیم زیرا هر کورس می تواند تعداد زیادی دسته بندی داشته باشد بنابراین برای هر دسته بندی یک ایندکس در الاستیک سرچ میسازیم، و در نهایت برای زبان زیر نویس ها هم یک نمایه گذاری انجام میدهیم یعنی هر زبان زیر نویس یک ایندکس از کورس ها در الاستیک سرچ دارد. به نظر میرسد که مقداری داپلیکیت دیتا در ایندکس ها وجود داشته باشد، به عنوان مثال یک کورس میتواند در ۱۰ ایندکس باشد، اما مهمترین نکته این است که با این کار پرفورمنس بهتری برای دسترسی به نتایج به دست می آوریم.

پاکسازی داده

برای پاکسازی داده ها بعد از لود کردن دیتای اولیه با استفاده از پانداز ابتدا قیمت سرتیفیکیت هایی که مشخص نیستند را برابر میانگین قرار میدهیم، همچنین کاراکترهایی که برای ذخیره سازی در الاستیک سرچ قابل قبول نیستند را به آندر لاین تغییر میدهیم.

در نهایت با استفاده از کتابخانه الاستیک سرچ پایتون کورس ها را در الاستیک سرچ اینسرت میکنیم که حدود یک و نیم ساعت طول میکشد.

کوئری ها

کوئری های خواسته شده را با استفاده از پایتون انجام داده و نتایج را ذخیره میکنیم.

ساختار پروژه

- Indexer.py

برنامه ایندکس کننده دیتای تولید شده توسط خزشگر

- Query.py

برنامه اجرا کننده کوئری‌های خواسته شده با نام گذاری از q1 تا q9، میتونید خود کوئری‌ها و ایندکس‌هایی روی آن اجرا میشوند را در این فایل مشاهده کنید.

- Req.txt

ریکویرمنت‌های پروژه برای اجرا شدن.

- Config.py

کانفیگ‌های مربوط به پروژه شامل آدرس الاستیکسرچ.

- Data_new.py

آخرین دیتای ماین شده از خزشگر.

- crawler/courses_spyder.py

خزشگری که دیتای اولیه کورس‌ها را استخراج میکند.

- results/

- ?.json

ریسپانس مربوط به هر کوئری با شماره ؟

- ?.query.json

کوئری مربوط به هر سوال با شماره ؟

افراد گروه

یوسف سلملیان

برنا کلهر

فرناز بنائیان

بهار عالمی

افاقیا محمدی