

## Search System

### آرش محرابی

### زهرا کولایی زاده

#### دیتابیس:

برای این پروژه از دیتابیس MySQL استفاده شده است.  
جدول news:

Column	Type
id	int
title_tokens	json
body_tokens	json
title_text	varchar(255)
description	varchar(255)
url	varchar(255)
date	varchar(255)

title\_tokens و body\_tokens توکن های پیش پردازش شده هستند که به صورت json ذخیره می شوند.  
بقیه ستون ها بدون پیش پردازش ذخیره می شوند. ( برای نمایش هنگام retrieve کردن result )

#### جدول tf\_idf:

Column	Type
id	int
news_id	int
token	varchar(100)
tf_idf	varchar(45)

news\_id در حقیقت یک foreign key به id در جدول news است.  
در این جدول tf\_idf را برای هر زوج (news\_id, token) حساب و ذخیره می کنیم. پس باید زوج نامبرده UNIQUE باشد.  
tf\_idf یک عدد double است ولی ترجیح دادم به صورت varchar ذخیره کنم.

#### کد:

هنگام کد زنی سعی شده از اصول کد نویسی تمیز عمو باب استفاده شود بنابراین نام variable ها و تابع ها تا حد امکان گویا هستند و تابع ها معمولاً کمتر از 5 6 خط هستند و هر کدام یک کار انجام می دهند.

#### Crawler.py:

این فایل به صفحه news.urmia.ac.ir می رود و در هر صفحه همه خبر ها را پیدا می کند. Title و Body آن ها را پیش پردازش می کند و در نهایت در دیتابیس ذخیره می کند. بعد از اتمام هر صفحه به صفحه بعدی می رود و همین کار ها را تکرار می کند.

#### پیش پردازش:

برای پیش پردازش کردن ( نورمالایز کردن، حذف stopword ها و tokenize کردن ) از ماژولی که خودم در طول تمرین ها دولوپ کردم درست شده که در پوشه normalize قرار دارد.

:tf\_idf.py

کار کلی این فایل پر کردن جدول tf\_idf است. ابتدا کل خبر ها را بر می دارد و بعد df و idf را برای کل توکن های موجود در کل خبر ها حساب می کند، سپس tf\_idf را برای هر زوج داکيومنت، توکن حساب می کند.

:db\_mysql.py

این فایل یک کلاس DB\_MYSQL دارد و هر وقت که بخواهم از کد پایتون به دیتابیس وصل شوم از این کلاس یک object می سازم. این کلاس تابع هایی دارد که با آن ها خبر ها و tf\_idf ها را در جدولشان درج می کنم.

:preprocess\_query.py

این تابع را هر وقت بخواهم چیزی را در PHP پیش پردازش کنم صدا می زنم و متنی که قرار است پیش پردازش شود را به صورت argument به این مازول می دهم. نتیجه پیش پردازش شده آن متن است که به صورت json برگردانده می شود.

:Web Application

قسمت وب اپلیکیشن را با لاراول ساختم.

:Route

Domain	Method	URI	Name	Action	Middleware
	GET HEAD	/	home	Closure	web
	GET HEAD	api/user		Closure	api auth:api
	GET HEAD	news	news.index	App\Http\Controllers\NewsController@index	web
	POST	news	news.store	App\Http\Controllers\NewsController@store	web
	GET HEAD	news/create	news.create	App\Http\Controllers\NewsController@create	web
	GET HEAD	news/{news}	news.show	App\Http\Controllers\NewsController@show	web
	PUT PATCH	news/{news}	news.update	App\Http\Controllers\NewsController@update	web
	DELETE	news/{news}	news.destroy	App\Http\Controllers\NewsController@destroy	web
	GET HEAD	news/{news}/edit	news.edit	App\Http\Controllers\NewsController@edit	web
	GET HEAD	search	search	App\Http\Controllers\SearchController@search	web

:SearchController

query ای که کاربر وارد می کند به این کنترلر می آید. ابتدا کوئری توسط فایل preprocess\_query.py پیش پردازش می شود. سپس از جدول tf\_idf سطر هایی که شامل توکن های کوئری باشند برداشته می شود.

سپس امتیاز TF-IDF برای هر داکيومنت مرتبط حساب می شود. سپس داکيومنت ها ( news ) بر اساس امتیازشان مرتب می شوند و نتایج به result view می روند تا نمایش داده شوند.

:NewsController

در این کنترلر عملیات CRUD روی News بر اساس RestAPI کنترلر می شود. Index پنج خبر آخر را بر می گردند.

Store خبری که کاربر وارد کرده را validate می کند. سپس اگر مشکلی نداشت پیش پردازش می کند و در نهایت در جدول news ذخیره می کند.

Update خبری که کاربر وارد کرده را validate می کند. سپس اگر مشکلی نداشت پیش پردازش می کند و در نهایت به جای news قبلی ذخیره می کند.

Destroy بر اساس آی دی خبر را پیدا کرده و از جدول news حذف می کند.

Create, show, edit ها تنها view برای انجام این کار ها را بر می گردانند.