



دانشگاه صنعتی امیرکبیر

(پلی تکنیک تهران)

دانشکده مهندسی کامپیوتر و فناوری اطلاعات

گزارش کارآموزی (هفته پنجم)

محل کارآموزی: شرکت سامانه گستر سحاب پرداز

نام استاد کارآموزی

دکتر مسعود صبائی

نام دانشجو

امیرمحمد پیرحسینلو

۹۵۳۱۰۱۴

تابستان ۱۳۹۸

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

فهرست مطالب

۱	پیاده سازی پروژه
۲	۱-۱ مقدمه
۲	۲-۱ پکیج‌ها
۲	۱-۲-۱ crawler
۲	۲-۲-۱ index
۲	۳-۲-۱ storage
۳	۴-۲-۱ test
۳	۵-۲-۱ utils
۳	۶-۲-۱ language
۳	۳-۱ کلاس‌ها
۳	۱-۳-۱ Crawler
۳	۲-۳-۱ Fetcher
۴	۳-۳-۱ LruCache
۴	۴-۳-۱ Parser
۵	۵-۳-۱ Elastic
۵	۶-۳-۱ HBase
۵	۷-۳-۱ Constants
۵	۸-۳-۱ Statistics
۵	۹-۳-۱ Pair
۵	۴-۱ استثناها
۵	۵-۱ اجرای برنامه
۶	۶-۱ موارد باقی مانده
۶	۷-۱ نتیجه گیری
۷	واژه‌نامه‌ی انگلیسی به فارسی

فصل اول

پیاده سازی پروژه

۱-۱ مقدمه

در این هفته به پیاده سازی و برنامه نویسی موتور جستجو پرداختیم. اعضای تیم به برنامه نویسی بخش های مختلف پروژه پرداختند و در انتها با ترکیب کدها که در بستر **git** ذخیره شده بودند، برنامه ی واحدی که همان موتور جستجو است ساخته شد. این برنامه از دو بخش تشکیل شده است:

- خزنده: یک برنامه که با زبان جاوا نوشته شده است و وظیفه ی آن خزش (**crawl**) کردن صفحات وب و ذخیره اطلاعات آن است

و

- **search.sh**: یک اسکریپت است که از کاربر **query** گرفته و آن را در بین اسناد جستجو می کند و نتیجه را برمی گرداند. نتیجه شامل اسنادی است که شامل **query** هستند یا عبارات **query** در آن ها حضور دارد. نتایج با توجه به امتیازشان به صورت نزولی مرتب شده و نمایش داده می شوند (بهترین سند یافت شده اول نمایش داده می شود).

در ادامه به توضیح پکیج ها و کلاس های پروژه می پردازیم.

۲-۱ پکیج ها

پروژه از پکیج های زیر تشکیل شده است:

۱-۲-۱ **crawler**

این پکیج شامل تابع **main** است. در این تابع صف لینک ها با لینک (^۱**URL**) های اولیه پر می شود و نخ (**thread**) آمارگیری (کلاس **Statistics**) اجرا می شود. همچنین به تعداد نیاز از کلاس های **Fetcher** و **Parser** که خود نخ هستند نمونه ساخته شده و اجرا می شوند.

این پکیج شامل کلاس های زیر است: (کلاس ها در ادامه توضیح داده خواهند شد).

• **Crawler**

• **Fetcher**

• **Parser**

• **LruCache**

۲-۲-۱ **index**

این پکیج شامل کلاس **Elastic** است که وظیفه برقراری ارتباط با **ElasticSearch** و تبادل داده با آن را برعهده دارد.

۳-۲-۱ **storage**

این پکیج شامل کلاس **HBase** است که وظیفه برقراری ارتباط با **HBase** و تبادل داده با آن را برعهده دارد.

^۱Uniform Resource Locator

test ۴-۲-۱

شامل کلاس‌هایی است که برای تست سایر کلاس‌های نرم‌افزار به کار می‌روند.

utils ۵-۲-۱

کلاس‌های زیر در این پکیج قرار می‌گیرند:

• **Constant**

• **Pair**

• **Prints**

• **Statistics**

این کلاس‌ها در ادامه توضیح داده می‌شوند.

language ۶-۲-۱

شامل کلاس‌های مورد نیاز برای تشخیص زبان یک متن است.

۳-۱ کلاس‌ها

برای پیاده سازی پروژه از کلاس‌های مختلفی استفاده شده است که در ادامه به تشریح آن‌ها می‌پردازیم:

Crawler ۱-۳-۱

تابع **main** در این کلاس قرار دارد. عملیات انجام شده در این تابع در قسمت پکیج **crawler** توضیح داده شد.

Fetcher ۲-۳-۱

این کلاس خود یک نخ است و وظیفه دارد که به طور متناوب از صف لینک‌ها (**URLs queue**)، لینک (**URL**) دریافت کرده و محتوای آن را دانلود کند سپس لینک و محتوایش (**document**) را در صف لینک-سند (**url-document**) قرار دهد. برای انجام این عملیات از توابع زیر استفاده شده است:

fetchURL

یک لینک از سر صف لینک‌ها برداشته و برمی‌گرداند.

getDomainIfLruAllowed

بررسی می‌کند که آیا به **Host** لینک مورد نظر در ۳۰ ثانیه اخیر درخواستی داده شده است یا خیر. در صورت مثبت بودن، این لینک مجدداً در انتهای صف لینک‌ها قرار می‌گیرد و لینک دیگری اخذ می‌شود.

fetch

محتوای یک لینک را دانلود کرده و برمی‌گرداند.

putFetchedData

لینک و محتوای آن را در صف لینک-سند قرار می‌دهد.

LruCache ۳-۳-۱

در این کلاس پیاده سازی یک **LruCache** انجام شده است.

Parser ۴-۳-۱

این کلاس خود یک نخ است که به طور متناوب از صف لینک-سند داده (لینک و سند مربوط به آن) را استخراج می‌کند. محتوای لینک را **parse** کرده و موارد زیر را از آن استخراج می‌کند:

- عنوان لینک (**title**)

- **anchor link** ها و **anchor text** ها

- تگ‌های **<p>** (پاراگراف‌ها)

پس از استخراج اطلاعات بالا، عملیات‌های زیر انجام می‌شود:

- قرار دادن لینک و عنوان آن و پاراگراف‌ها در **ElasticSearch**

- قرار دادن لینک، **anchor link** ها و **anchor text** ها در **HBase**

- تحویل **anchor link** ها به **Kafka**

برای انجام عملیات‌های بالا از توابع زیر استفاده شده است:

takeFetchedData

از سر صف لینک-سند داده برداشته و برمی‌گرداند.

extractLinkAnchors

anchor link های یک سند را استخراج می‌کند و برمی‌گرداند.

putToElastic

لینک و عنوان آن و پاراگراف‌های سند مربوط به لینک را در **ElasticSearch** قرار می‌دهد.

putAnchorsToHBase

لینک، **anchor link** های درون سند مربوط به لینک و **anchor text** ها را در **HBase** قرار می‌دهد.

putToKafka

anchor link های مربوط به یک سند را به **Kafka** تحویل می‌دهد.

۵-۳-۱ Elastic

رابط برنامه و موتور جستجوی **ElasticSearch** است. مهم ترین تابع آن **indexData** است که لینک و عنوان آن و پاراگراف های سند مربوط به آن را دریافت کرده و داخل **ElasticSearch** قرار می دهد.

۶-۳-۱ HBase

رابط برنامه و پایگاه داده ی **HBase** است. مهم ترین تابع آن **insertLinks** است که لینک، **anchor link** های درون سند مربوط به لینک و **anchor text** ها را گرفته و داخل **HBase** قرار می دهد.

۷-۳-۱ Constants

مقادیر ثابت و پارامترهای برنامه در این کلاس قرار دارند.

۸-۳-۱ Statistics

کلاسی است که خود یک نخ است و وظیفه دارد به صورت تناوبی وضعیت سایر نخ ها (**Parser** ها و **Fetcher** ها) را گزارش دهد.

۹-۳-۱ Pair

پیاده سازی یک ساختار داده برای نگهداری یک کلید و مقدار متناظر آن است.

۴-۱ استثناها

موتور جستجوی ساخته شده تا حد بسیار زیادی مطابق با اطلاعات داده شده در این گزارش و گزارش های قبلی است و فقط در چند مورد زیر تفاوت دارد که علت آن نیز ذکر شده است:

- عدم استفاده از پکیج تشخیص زبان: به دلیل خطای بالای تشخیص زبان و زمان زیادی هم که می گرفت، از کلاس های این پکیج استفاده نشد. در اسرع وقت راه حلی برای مشکل پیش آمده پیدا می کنیم.
- عدم استفاده از **Kafka**: به جای **Kafka** ، **anchor link** ها را مستقیماً در صف لینک ها قرار دادیم و عدم مشاهده آن ها تاکنون را به یک صف دیگر که مدیریت آن بر عهده **HBase** است، سپردیم. دلیل عدم استفاده از **Kafka** ، عدم پیگیربندی مناسب و از کار افتادن بی مورد آن بود. در اسرع وقت سعی می کنیم این مشکل را حل کنیم.

۵-۱ اجرای برنامه

برنامه توسط ابزار مدیریت پکیج ^۲**maven** ساخته (build) شده است. نسخه ای که در کنار گزارش قرار دارد نسخه ای است که بر روی یک سرور اجرا می شود. برای مشاهده نسخه اصلی کد که در دو سرور اجرا می شود به آدرس زیر مراجعه کنید:

<https://github.com/Ahmad535353/NimboSearchEngine>

برای اجرای بخش خزنده برنامه کافی است روند زیر طی شود:

- اجرای **Hadoop** ، **HBase** و **ElasticSearch**

^۲<https://maven.apache.org>

- ایجاد یک جدول به نام CrawlerTable در HBase با فامیلی ستون ^۳ data

- java -jar Crawler.jar

برای ارسال query و دریافت پاسخ، کافی است فایل اجرایی **search.sh** را اجرا کنید. با اجرای فایل از شما درخواست می شود تا query را وارد کنید. پس از قرار دادن query و زدن دکمه **Enter**، پاسخ دریافت و نمایش داده می شود.

۶-۱ موارد باقی مانده

به دلیل ضیق وقت، فرصت نشد تا موارد زیر پیاده سازی شوند. سعی می شود در اسرع وقت این موارد پیاده سازی شوند:

- پشتیبانی گرفتن از صف های لینک، لینک-سند و صف کلاس HBase و جلوگیری از دست رفتن داده ها (data loss) هنگام قطع شدن برق.
- فیلتر کردن برخی پروتکل ها مانند FTP و ...
- برطرف کردن مشکل Kafka
- پیاده سازی و یا پیدا کردن یک ماژول تشخیص دهنده زبان یک متن
- جایگزین کردن کلاس های منسوخ شده (deprecated) برنامه مانند کلاس های استفاده شده در کلاس Elastic و ...

۷-۱ نتیجه گیری

با اتمام پیاده سازی، یک موتور جستجو در اختیار ماست که فقط «کار» می کند. نتیجه query ها اصلا خوب نیستند. تیم شرکت به ما توصیه کردند که از الگوریتم ^۴Page Rank و مدل برنامه نویسی MapReduce برای امتیازدهی سایت ها و رتبه بندی آن ها استفاده کنیم و همچنین نواقص قبلی را برطرف کنیم. برای پیاده سازی Page Rank از فریم ورک ^۵Spark استفاده می کنیم. پیاده سازی باقی مانده و رفع نواقص قبلی در هفته آینده انجام خواهد شد.

^۳Column Familiy

^۴<https://en.wikipedia.org/wiki/PageRank>

^۵<https://spark.apache.org/>

واژه‌نامه‌ی انگلیسی به فارسی

B	deprecated منسوخ شده
build ساختن. در این جا منظور ساخت فایل jar است.	document سند
C	T
crawl خزیدن	thread نخ
D	title عنوان