

Elasticsearch

Elasticsearch یک موتور جست و جو متن باز (Open Source) و بلادرنگ (real-time) برای تحلیل متون، جستجو و ذخیره سازی انواع داده های نیمه ساختاریافته به صورت توزیع پذیر است .

موتورهای جستجو

موتورهای جستجو (search engines) ابزارهایی هستند که به طور خاص برای نیازمندی‌های پیچیده‌ی کاربران در حوزه‌ی جستجو طراحی شده‌اند.

در این بین مواردی هستند که به آنها online search engine یا موتور جستجوی آنلاین گفته می‌شود که به طور خاص روی محتوای وب و فعالیت‌های آنلاین تمرکز داشته و داده‌ی ورودی آنها محتواهایی است که افراد مختلف در بستر شبکه‌ی اینترنت منتشر می‌کنند.

روش کار موتورهای جستجو

فرایند کار موتورهای جستجو شامل ۳ مرحله‌ی اساسی است:

۱) جمع‌آوری داده:

این مرحله در موتورهای جستجوی آنلاین به مرحله‌ی **Crawling** معروف است. معنای لغوی crawl خزیدن است و به مرحله‌ای گفته می‌شود که طی آن موتورهای جستجو با استفاده از ابزارهای ربات‌هایی که در اختیار دارند، شروع به **خواندن** اطلاعات صفحه‌های وبسایت‌های مختلف کرده و داده‌های منتشر شده در آنها را **جمع‌آوری می‌کنند**. برخی موتورهای جستجو که به صورت آنلاین کار نمی‌کنند، در این مرحله به صورت **دستی (manually)** داده‌های ورودی آنها **تامین** می‌شود.

۲) شاخص‌گذاری :

این مرحله در اصطلاح (indexing) شاخص‌گذاری نامیده می‌شود.

کلمات کلیدی داده‌ها مورد بررسی قرار گرفته و برای ثبت شدن در پایگاه داده دسته‌بندی و استخراج می‌شوند.

۳) رتبه بندی:

در زمان جستجوی یک عبارت، مجموعه ای از الگوریتمها استفاده می شود تا امتیازی برای میزان ارتباط هر یک از داده ها با عبارت جستجو شده به دست آید. در نهایت بر مبنای امتیاز به دست آمده نتایج مرتب شده و مرتب ترین آنها به کاربر نمایش داده می شود.

Elasticsearch یک موتور جست و جو متن باز (Open Source) و بلادرنگ (real-time) برای تحلیل متون، جستجو و ذخیره سازی انواع داده های نیمه ساختاریافته به صورت توزیع پذیر (distributable) است.

متن باز بودن (open source) : Elasticsearch با زبان برنامه نویسی جاوا توسعه داده شده است و سورس کد آن در سایت github.com قابل دسترس است.

بلادرنگ بودن : Elasticsearch پس از ذخیره داده ها در کمترین فاصله آنها را برای جستجو آماده می کند. همان طور که در بخش قبل توضیح داده شد، دومین مرحله ی کار موتورهای جستجو، indexing یا شاخص گذاری کلمات کلیدی آنهاست.

Elasticsearch به طور پیش فرض یک ثانیه پس از ثبت داده ها شروع به index کردن آنها می کند.

Elasticsearch یک موتور جست و جو متن باز (Open Source) و بلادرنگ (real-time) برای تحلیل متون، جستجو و ذخیره‌سازی انواع داده های نیمه ساختاریافته به صورت توزیع پذیر (distributable) است.

تحلیل و جستجوی : منظور از تحلیل کردن داده ها، تبدیل آنها به فرمتی بهینه برای عملیات full text search است. پس از آن Elasticsearch می تواند به سرعت تمامی داده های مرتبط با یک عبارت مورد نظر را شناسایی و امتیازدهی کند.

پشتیبانی از داده های نیمه ساختاریافته : منظور از داده ی نیمه ساختاریافته این است که تمامی داده ها لزومی ندارد ساختار کاملا یکسان داشته باشند. برای مثال فرض کنید داده های log بازدید را داشته باشیم و برخی از این log ها دارای یک فیلد به نام event_id باشند و برخی دیگر آن را نداشته باشند اما در عین حال همگی آنها شامل فیلدهای مشترک source، created_at و type باشند. در داده های ساختاریافته (structured) که به طور معمول در پایگاه داده ی رابطه ای از آنها استفاده می شود، تمامی داده ها باید ساختار کاملا یکسانی داشته باشند و برعکس در داده های بدون ساختار (unstructured) هیچ شمای مشخصی در هیچ یک از فیلدهای داده ای آن وجود ندارد مانند داده های ویدیویی.

Elasticsearch یک موتور جست و جو متن باز (Open Source) و بلادرنگ (real-time) برای تحلیل متون، جستجو و ذخیره‌سازی انواع داده‌های نیمه ساختاریافته به صورت توزیع پذیر (distributable) است.

توزیع پذیری : Elasticsearch به گونه‌ای طراحی شده است که بتوان منابع مورد استفاده‌ی آن را با کمترین دردسر گسترش داد تا در صورت نیاز بتواند از حجم بالاتر داده‌ها و پردازش‌های سنگین پشتیبانی کند. طراحی Elasticsearch بر اساس معماری cluster یا خوشه است. در این معماری تعدادی سرور که در یک شبکه به یکدیگر متصل هستند به عنوان یک سیستم واحد عمل کرده‌اند و ضمن ارتباط با یکدیگر درخواست‌های دریافت شده را پاسخ می‌دهند

ساختار ذخیره و بازیابی داده ها در Elasticsearch

Elasticsearch داده ها را در فرمت JSON می شناسد. در Elasticsearch به جای ساختار کلاسیک جدول ها و ستون ها در پایگاه داده ی رابطه ای، مفاهیمی به نام **index** و **document** معرفی می شود.

می توان **index** را به عنوان جدول و **document** را ردیفی از جدول تصور کرد.

هر **Index مجموعه ای** از **document** ها است و داده هایی که در یک **document** ذخیره می شوند به صورت نیمه ساختاریافته هستند.

زمانیکه یک **document** در Elasticsearch ثبت می شود، علاوه بر اصل (**source**) داده های ورودی، بخشی از آن توسط تحلیلگرها

(**Analyzers**) مورد پردازش قرار می گیرد و پس از پردازش به اجزایی کوچک تر (**terms**) تبدیل می شود. هدف از پردازش توسط

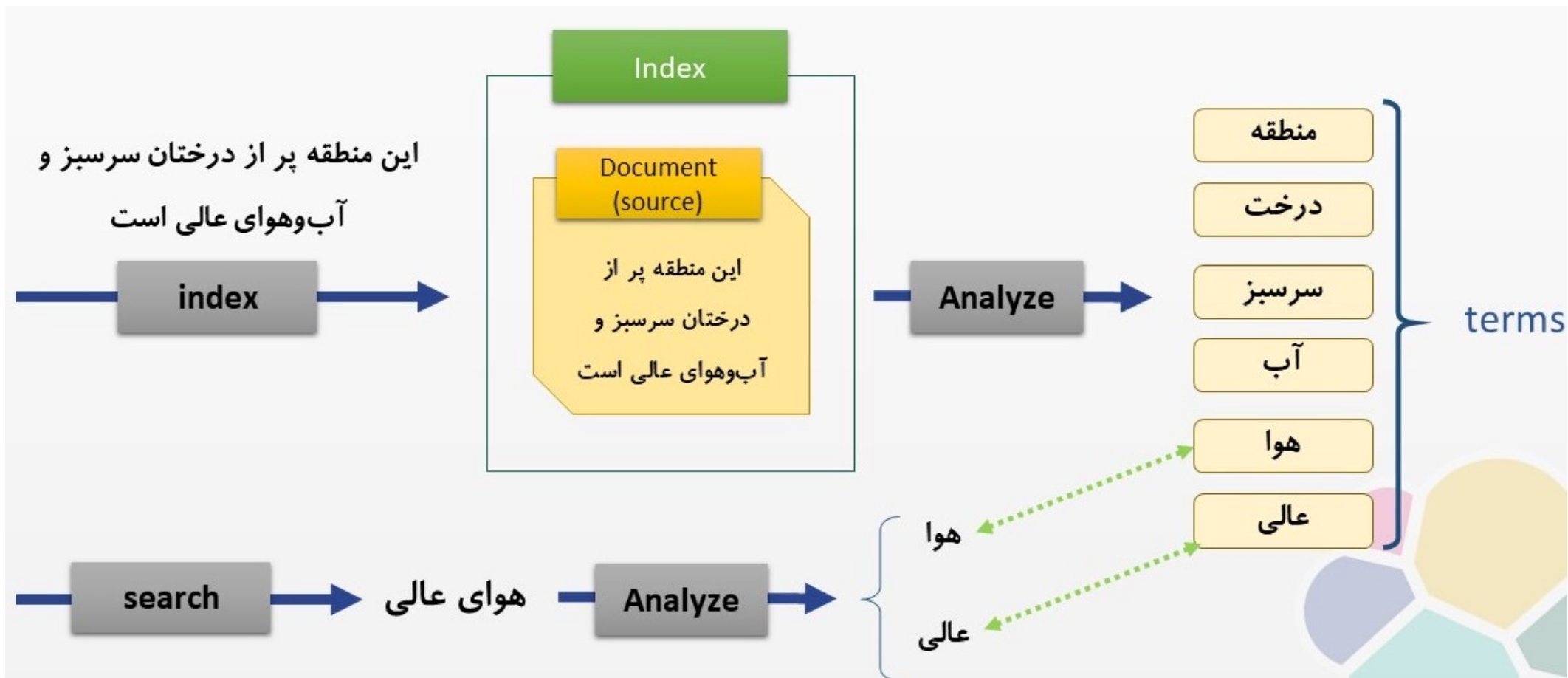
تحلیلگرها، آماده سازی داده ها برای شرکت در عملیات جستجو است. شاید بتوان گفت بدون تحلیلگرها، Elasticsearch هیچ حرفی

برای گفتن در زمینه ی جستجو نداشت!

پس از ثبت داده ها، نوبت به جستجوی آنها می‌رسد. ES قدرت خود در جستجو را، از کتابخانه‌ی **Apache Lucene Search** می‌گیرد. API های **Elasticsearch** به سادگی به شما اجازه می‌دهند تا **Query** های جستجو را به **Elasticsearch** ارسال کنید.

همچنین تیم **Elasticsearch** اقدام به انتشار **Client** های مبتنی بر زبان های برنامه نویسی مختلف کرده‌اند تا به سادگی بتوان در نرم افزارهای گوناگون از API های ارائه شده‌ی **ES** استفاده کرد.

نمای کلی از فرایند تجزیه و تحلیل داده ها



طراحی مبتنی بر اصول "همیشه در دسترس بودن" و "توزیع پذیری"

یکی از نکته های مهم در Elasticsearch در نظر گرفتن طراحی توزیع پذیر (Distributable) آن است به طوری که امکان توسعه ی آن را بسیار ساده و پشتیبانی از حجم بالای داده ها و درخواست ها را امکان پذیر می کند.

افزون بر توزیع پذیری، ساختار Elasticsearch به گونه ای طراحی شده است که در صورت از دست رفتن قسمتی از منابع مورد استفاده ی ES، دسترسی به داده ها تا حد امکان همچنان برقرار باشد. به عبارتی همواره نسخه ی کپی از داده ها به صورت توزیع شده در منابع (سرورهای) مختلف وجود داشته باشد تا بتوان از آن ها در مواقع نیاز استفاده کرد.

Elasticsearch بر مبنای معماری فوشه بندی (cluster) طراحی شده است. در این معماری مجموعه ای از سرورها که گره نامیده می شوند، در یک شبکه ی داخلی به یکدیگر متصل شده و ضمن ارتباط درونی با یکدیگر، درخواست هایی که به سمت آنها ارسال می شود را، به گونه ای پاسخ می دهند که کاربر تصور می کند درخواست ارسال شده با یک سیستم یکپارچه پاسخ داده می شود. به عبارتی در این معماری فرقی نمی کند درخواست توسط کدام گره یا سرور دریافت شود و در صورت لزوم گره ها اطلاعاتی را مابین یکدیگر تبادل فوهند کرد و در نهایت پاسخ کاربر داده فوهند شد. بنابراین در Elasticsearch شما به سادگی می توانید در صورت نیاز، گرهی جدیدی به فوشه اضافه کنید و Elasticsearch به طور خودکار بار درخواست ها و داده ها را در سطح فوشه تعدیل می کند.

برای فهم بهتر مفاهیم فرض کنید لیستی از اطلاعات پروازهای داخلی شرکت هواپیمایی ماهان را به عنوان داده ی نمونه داشته باشیم تا در تعاریف زیر از این داده ی نمونه استفاده کنیم:

Document: یک واحد اطلاعات به فرمت JSON است که در Elasticsearch ذخیره می شود. در documentها اصل (source) داده ها بدون هیچ تغییری ذخیره می شود. در مثال داده ی نمونه، هر document برابر اطلاعات یکی از پروازها خواهد بود.

Field: فیلدهای اطلاعاتی یا همان data field ها که به صورت key و value در document ذخیره می شوند. برای مثال، ساعت پرواز یا مبدا و مقصد پرواز مثالی هایی از field هستند.

Index: به مجموعه ای از document های مربوط به یک موضوع مشخص، index گفته می شود. برای مثال مجموعه ی اطلاعات همه ی پروازها یک index خواهد بود.

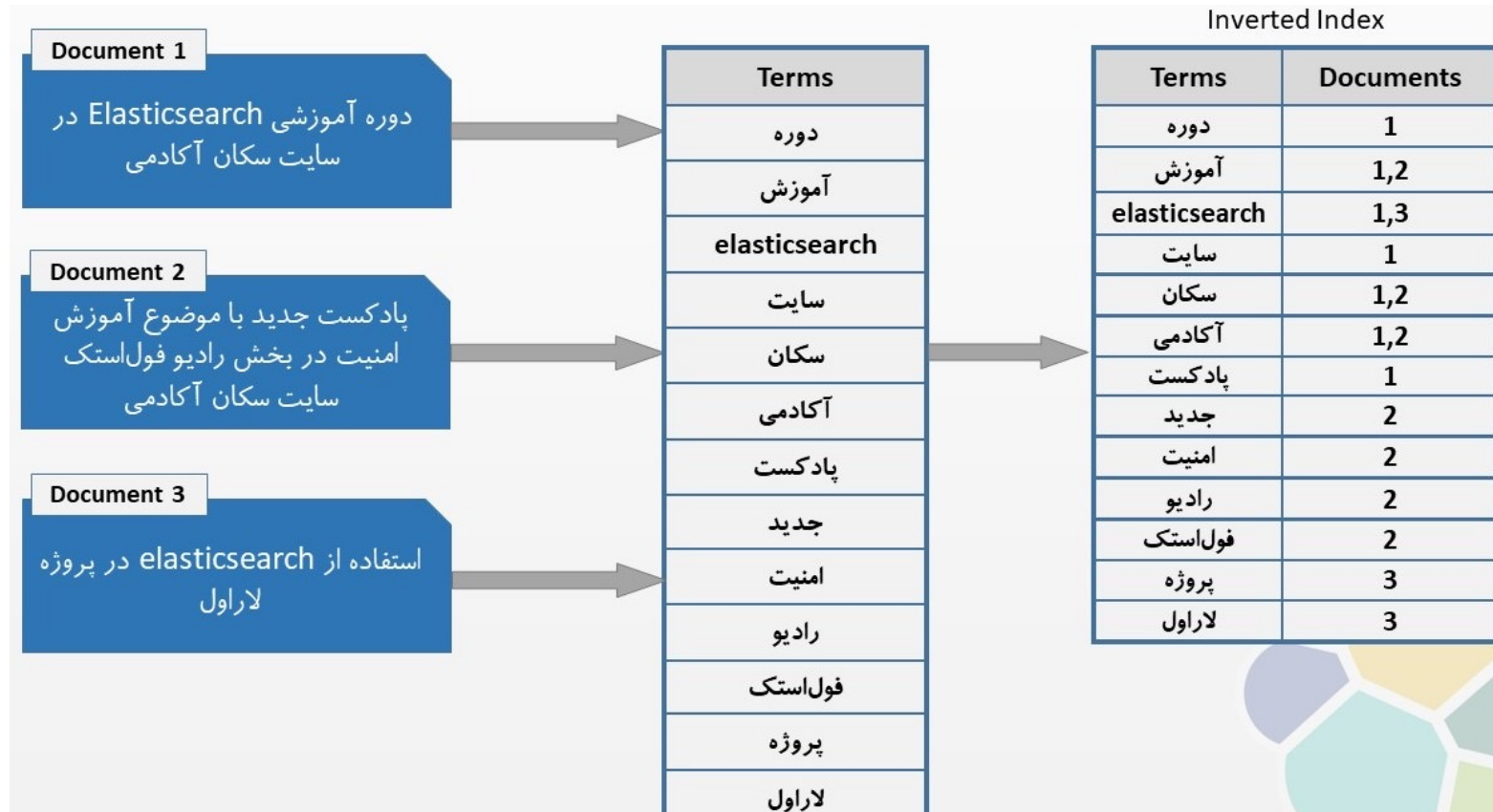
Indexing: عمل ثبت یک document در index را اصطلاحاً indexing یا شافص گذاری می گویند. طی این فرایند علاوه بر اصل داده های document، کلمات کلیدی متون شناسایی و ذخیره می شوند.

Mapping: تعریف ساختار داده ها در فیلدهای اطلاعاتی document های یک index را mapping می گویند که شامل مواردی همچون نوع داده، نوع تملیلگر و سایر اطلاعات اضافی در مورد هر یک از field ها می شود. Mapping را می توان مشابه schema در پایگاه داده های رابطه ای در نظر گرفت. در نمونه ی داده های پرواز، اینکه برای مثال فیلد ساعت پرواز از نوع زمان، شماره پرواز به عنوان کلید یکتا و مبدا و مقصد پرواز از نوع کلمات کلیدی باشند، در واقع mapping ایندکس پروازها را مشخص می کند.

Token (Term): هنگامی که یک متن (Text،) توسط تملیلگر ها در Elasticsearch مورد پردازش قرار گیرد، خروجی آن مجموعه ای از term ها است که معمولا (نه متما همیشه!) کلمات مجزا از هم و معنادار بوده و برای شرکت در فرایند جستجو بسیار بهینه هستند. برای مثال عبارت "من یک برنامه نویس حرفه ای هستم" ممکن است تبدیل به term های { من، برنامه، نویس، حرفه } شود.

Cluster: یک یا مجموعه ای از node های Elasticsearch که با یکدیگر در ارتباط باشند، یک cluster یا خوشه را تشکیل می دهد. node ها در سطح یک cluster با یکدیگر در ارتباط بوده و یکدیگر را می شناسند.

Inverted Index : یک ساختار داده‌ی بهینه شده در حافظه ره است که طبق آن به سرعت مشخص می‌شود هر term یا کلمه کلیدی در کدام documentها تکرار شده است. اگر Index را همانند فهرست موضوعی ابتدای کتابها در نظر بگیریم، Inverted Index معادل فهرست واژگان انتهای کتاب است که جستجوی سریع یک واژه بر مبنای اینکه در چه صفحه‌هایی تکرار شده است را ممکن می‌کند!



شمای کلی از ساختار یک خوشه فرضی

- در این خوشه شرایط زیر برقرار است:
- شامل ۲ گره (سرور) به نام‌های node_1 و node_2 می‌باشد.
- شامل دو index به نام‌های index_1 و index_2 و هر یک از index ها شامل دو shard می‌باشد.
- برای index_1 هیچگونه replica shard یا نسخه‌ی کپی در نظر گرفته نشده است و برای index_2 تعیین شده است تا یک نسخه‌ی پشتیبان از shard های اصلی داشته باشد.

