

#Elasticsearch  
Erick Madrigal Zavala  
2018146983

## Datos en: documentos e índices

Elasticsearch es un almacenamiento de documentos distribuidos. En lugar de almacenar la información en forma de tablas, Elasticsearch almacena complejas estructuras de datos que se serializan como documentos JSON. Elasticsearch usa una estructura de datos llamada índices invertidos que soporta búsquedas de texto de manera muy rápida. Un índice inverso lista cada palabra única que aparece en un documento e identifica cada documento en que la palabra se encuentra.

Elasticsearch indexa todos los datos de cada campo y cada campo de indexado tiene una estructura de datos dedicada y optimizada. Elasticsearch también tiene la habilidad de ser "schema-less", que significa que los documentos pueden ser indexados sin tener que decir explícitamente como manejar cada uno de los campos que posee el documento. Cuando se activa el mapeo dinámico, Elasticsearch detecta y añade campos al índice automáticamente. El usuario puede definir reglas de control para el mapeo dinámico y además tomar el control total de como se guardan los datos en los índices por medio de definir explícitamente el mapeo.

Definir sus propios mapeos le permite:

- \* Distinguir entre string con el texto completo y campos que contengan un valor exacto de un string
- \* Realizar un análisis de lenguaje específico del texto.
- \* Optimizar campos para coincidencias parciales.
- \* Usar formatos de fechas personalizadas.
- \* Usar tipos de datos como "geo\_points" y "geo\_shapes" las cuales no podrían ser detectados automáticamente.

El análisis de cadenas que se aplica a un campo con todo el texto durante la indexación es también usado durante el tiempo de búsqueda. Cuando se realiza un query de un campo de texto, el query del texto pasa por el mismo análisis antes de que los términos sean buscados en el índice

## Salida de información: búsqueda y análisis

El REST API de Elasticsearch soporta queries estructurados queries de "full-text" y queries complejos que combinan los dos. Los queries estructurados son similares a los queries que se construyen para SQL. Los queries de "full-text" buscan todos los documentos que concuerden con el string del query y los devuelven ordenados por relevancia. Además soporta búsqueda de términos individuales y da sugerencias para el autocompletado.

El usuario además puede construir queries con el estilo SQL para buscar y agregar datos de manera nativa a Elasticsearch, los drivers de JDBC y ODBC permiten un amplio rango de aplicaciones de terceros que pueden interactuar con Elasticsearch vía SQL.

## Analizando los datos

Las agregaciones de Elasticsearch le permiten al usuario crear resúmenes completos de los datos y obtener información sobre métricas, patrones y tendencias clave.

Debido a que las agregaciones aprovechan las mismas estructuras de datos que las búsquedas, también son muy rápidas. Esto le permite analizar y visualizar sus datos en tiempo

real. Los informes y cuadros de mando se actualizan a medida que cambian los datos para poder tomar medidas basadas en la información más reciente. Además, las agregaciones funcionan junto con las solicitudes de búsqueda. Puede buscar documentos, filtrar resultados y realizar análisis al mismo tiempo, sobre los mismos datos, en una única solicitud. Puede utilizar funciones de "Machine Learning" para crear líneas de base precisas del comportamiento normal de sus datos e identificar patrones anómalos. Con el "Machine Learning", puede detectar:

- \* Anomalías relacionadas con desviaciones temporales en valores, recuentos o frecuencias.
- \* Rareza estadística.
- \* Comportamientos inusuales para un miembro de una población.

## Escalabilidad y Resiliencia: clusters, nodos y shards

Se pueden añadir servidores (nodos) a un clúster para aumentar su capacidad y Elasticsearch distribuye automáticamente los datos y la carga de consultas entre todos los nodos disponibles. En realidad, un índice de Elasticsearch no es más que una agrupación lógica de uno o más shards físicos, donde cada shard es en realidad un índice autónomo. Al distribuir los documentos de un índice en varios shards y distribuir esos shards en varios nodos, Elasticsearch puede garantizar la redundancia, que protege contra fallos de hardware y aumenta la capacidad de consulta a medida que se añaden nodos a un clúster. A medida que el clúster crece (o decrece), Elasticsearch migra automáticamente los fragmentos para reequilibrar el clúster.

Existen dos tipos de shards: primarios y réplicas. Cada documento de un índice pertenece a un shard primario. Un shard réplica es una copia de un shard primario. Las réplicas proporcionan copias redundantes de sus datos para protegerlos contra fallos de hardware y aumentar la capacidad de servir peticiones de lectura como la búsqueda o recuperación de un documento. El número de shards réplica puede modificarse en cualquier momento, sin interrumpir las operaciones de indexación o consulta.

Hay una serie de consideraciones de rendimiento y compensaciones con respecto al tamaño de los fragmentos y al número de fragmentos primarios configurados para un índice. Consultar muchos fragmentos pequeños hace que el procesamiento por fragmento sea más rápido, pero más consultas significa más sobrecarga, por lo que consultar un menor número de fragmentos más grandes podría ser más rápido.

Como punto de partida:

- \* Intenta mantener el tamaño medio de los fragmentos entre unos pocos GB y unas decenas de GB. Para casos de uso con datos basados en el tiempo, es habitual ver fragmentos de entre 20 GB y 40 GB.
- \* Evita el problema de los gazillion shards. El número de fragmentos que puede contener un nodo es proporcional al espacio disponible en el montón. Como regla general, el número de shards por GB de espacio de heap debe ser inferior a 20. Los nodos de un clúster necesitan conexiones buenas y fiables entre sí. Para ofrecer mejores conexiones, los nodos suelen ubicarse en el mismo centro de datos o en centros de datos cercanos. Sin embargo, para mantener una alta disponibilidad, también es necesario evitar cualquier punto único de fallo. En caso de que se produzca un corte importante en una ubicación, los servidores de otra ubicación deben ser capaces de tomar el relevo. ¿Cuál es la solución? La replicación entre clústeres (CCR).