

Date _____ 20 _____

Section 2: mapping & indexing data

~~Mapping~~ schema definition

field types

field index do you want this field indexed for full text
search
analyzed / not analyzed

~~field analyzers~~

define tokenizers & token filter

standard whitespace etc

(how to split strings)

lowercase
uppercase etc

~~Analyzers~~

character filters

remove html encoding etc.

Tokenizer

split strings on whitespace/punctuation/
non letters etc.

Tokenfilter

lowercasing, stemming, synonym, stopwords

choices for Analyzers

standard

split on word boundaries, remove punctuation,
lowercases, good choice if language is unknown

simple splits on anything that isn't a letter & lowercase

whitespace splits on whitespace but doesn't lowercase

language accounts for language specific stopwords
& stemming

Date 20

only need to define mapping if something need to
be specifically set to a certain type
or else type inference will be used

Versions

Every document has a version field

Elasticsearch documents are immutable

When you update an existing document:

- new doc created w/ incremented _version
- . old doc marked for deletion

Date

20

Dealing with Concurrency

- Optimistic concurrency control

- sequence number

- primary shard that owns sequence

use retry-on-conflicts=N to auto retry

Controlling full text search

Using Analyzers

- sometimes text fields should be exact-match

- use keyword mapping instead of text

- search on analyzed text fields will return anything remotely relevant

- depending on the analyzer results will be case insensitive, stemmed, stop words removed, synonyms applied etc.

- Searches w/ multi field terms need not match them all

Date _____ 20 _____

Data Modeling

Strategies for relational data

Normalized data

Look up rating
(moviesLens dataset)

Rating
-userId
-movieId
-rating

Look up → title

Movie
-movieId
-title
-genres

- Min storage, makes it easy to change
- require 2 queries, storage cheap

doubles traffic on cluster - page response time

Denormalized data

look up rating →

Rating
userId
rating
title

duplicated titles, 1 query
changing title diffcult

Date _____ 20 _____

Parent / Child relationship

Movie franchises etc

Flattened Datatype

scaling w/ large # of subfields

flattened subfield: avoid every sub field ~~as individual~~ as individual field use one field subfield

elastic search's performance suffers if too many inner fields - mapping explosion

use when unknown or large # of fields

cluster collection of elastic search nodes
cluster state passed b/w ~~nodes~~ so clusters nodes

run smoothly. master node sends node state to other nodes, they ack

Date _____ 20 _____

after each cluster state change nodes need
to be synced

frequently adding new fields to index
causes:

- cluster state to grow
- triggers cluster state updates across
all nodes
which can cause delays
~~attempts~~

w/o update cluster state nodes cant index, search
etc can cause

- memory issues
- poor performance
- cluster crash (mapping explosion)

map entire object w/ inner field into
single field \rightarrow flattened

Date 20

limitation of flattened

Keys treated as keywords
no analyzers and tokenizer

Supported queries for flattened data types

- term, terms & term-set
- prefix
- range (non numerical range ops)
- match & multi-match (supply exact keywords)
- query_string & simple_query_string
- exist

result highlighting feature not enabled
for flattened data type

Date _____ 20 _____

Mapping Exceptions

~~Mapping~~



Defining how JSON doc will be stored

Res:

actual metadata resulting from the definition process

Explicit Mapping

predefined fields & types

Dynamic Mapping

fields & types defined by ES

Date _____ 20 _____

What could go wrong?

Explicit

Mapping Exception when mismatch

Dynamic

mapping explosion

Dead Queue Pattern

store failed documents
in separate queue

that throw
mapping
exception

handle on app level or use logstash DL queue

allows to still process
failed docs

Date _____ 20 _____

limits of mapping

Default # 1000 # of fields in a mapping

changing to greater than 1000 may
cause performance degradation & high
memory pressure