

HBase Table layout (v0.2)

Grouperfish clustering engine



documents

	row key	utf8	$f_1(\text{namespace}, \text{collection_key}, \text{id})^*$
main	namespace	utf8	namespace
	collection_key	utf8	name of the collection
	id	utf8	document id (given by the application)
	text	utf8	document text
	member_of	utf8	last assigned cluster label
processing***	id	utf8	<i>document id (given by the application), reduncant for scans</i>
	vector_idf	byte[]	<i>serialization of a mahout sparse feature vector (matching the collection dictionary)</i>
	vector_tfidf	byte[]	<i>corresponding tf/idf vector</i>

clusters

	row key	utf8	$f_2(\text{namespace}, \text{collection_key}, \$\text{conf}, \text{timestamp}, \text{label})^{**}$
	namespace	utf8	namespace
	collection_key	utf8	name of the collection
	timestamp	utf8	When this cluster was fully rebuilt. Matches the configuration:\$conf:rebuilt item of this conf.
	label	utf8	numeric id or (better) descriptive text label
	size	utf8	current size of the cluster
	\$documentid	utf8	the score for the given document in the cluster

collections

	row key	utf8	$f_3(\text{namespace}, \text{collection_key})$
main	namespace	utf8	namespace
	collection_key	utf8	name of the collection
	size	utf8	base 10 integer
	modified	utf8	base 10 unix timestamp / last document addition time
	configuration:\$conf:rebuilt	utf8	base10 unix ts / when this configuration was fully rebuilt
	configuration:\$conf:processed	utf8	base10 unix ts / when this configuration was updated
processing	dictionary	byte[]	<i>dictionary descibing the document's (TF) IDF vector features</i>
	index	byte[]	<i>inverted index (for similarity based algorithms)</i>

* The values of $f_1(\dots)$ for the same collection are all prefixed with the same $f'_1(\text{namespace}, \text{collection_key})$. This allows to quickly scan for all documents of a collection.

** The values of $f_2(\dots)$ for the same cluster are all prefixed with the same $f'_2(\text{namespace}, \text{collection_key}, \$\text{conf}, \text{ts})$. This allows to quickly get all documents in a cluster.

Note that for each collection we can store a clustering for each (conf, ts) combination. This allows to store a new clustering while serving the previous one.

*** The *processing* family is not used yet (only full rebuilds are supported).