

git_comments:

1. * Holds statistics for a DocValues field.
2. * * Called after #{@link DocValuesStats#accumulate(int)} was processed and verified that the document has a value for * the field. Implementations should update the statistics based on the value of the current document. * * @param count * the updated number of documents with value for this field.
3. * The number of documents which have a value of the field.
4. * The maximum value of the field. Undefined when #{@link #count} is zero.
5. * Returns whether the given document has a value for the requested DocValues field.
6. * The field for which these stats were computed.
7. * Licensed to the Apache Software Foundation (ASF) under one or more * contributor license agreements. See the NOTICE file distributed with * this work for additional information regarding copyright ownership. * The ASF licenses this file to You under the Apache License, Version 2.0 * (the "License"); you may not use this file except in compliance with * the License. You may obtain a copy of the License at * * <http://www.apache.org/licenses/LICENSE-2.0> * * Unless required by applicable law or agreed to in writing, software * distributed under the License is distributed on an "AS IS" BASIS, * WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. * See the License for the specific language governing permissions and * limitations under the License.
8. * The mean of all values of the field. Undefined when #{@link #count} is zero.
9. * Holds DocValues statistics for a numeric field storing {@code double} values.
10. * * Initializes this object with the given reader context. Returns whether stats can be computed for this segment (i.e. * it does have the requested DocValues field).
11. * Holds DocValues statistics for a numeric field storing {@code long} values.
12. * Holds statistics for a numeric DocValues field.
13. * The minimum value of the field. Undefined when #{@link #count} is zero.
14. **comment:** * The number of documents which do not have a value of the field.
label: documentation
15. All matching documents in this reader are missing a value
16. * A #{@link Collector} which computes statistics for a DocValues field.
17. * Licensed to the Apache Software Foundation (ASF) under one or more * contributor license agreements. See the NOTICE file distributed with * this work for additional information regarding copyright ownership. * The ASF licenses this file to You under the Apache License, Version 2.0 * (the "License"); you may not use this file except in compliance with * the License. You may obtain a copy of the License at * * <http://www.apache.org/licenses/LICENSE-2.0> * * Unless required by applicable law or agreed to in writing, software * distributed under the License is distributed on an "AS IS" BASIS, * WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. * See the License for the specific language governing permissions and * limitations under the License.
18. * Creates a collector to compute statistics for a DocValues field using the given {@code stats}.
19. Stats cannot be computed for this segment, therefore consider all matching documents as a 'miss'.
20. * Licensed to the Apache Software Foundation (ASF) under one or more * contributor license agreements. See the NOTICE file distributed with * this work for additional information regarding copyright ownership. * The ASF licenses this file to You under the Apache License, Version 2.0 * (the "License"); you may not use this file except in compliance with * the License. You may obtain a copy of the License at * * <http://www.apache.org/licenses/LICENSE-2.0> * * Unless required by applicable law or agreed to in writing, software * distributed under the License is distributed on an "AS IS" BASIS, * WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. * See the License for the specific language governing permissions and * limitations under the License.
21. 20% of cases delete some docs
22. * Unit tests for #{@link DocValuesStatsCollector}.
23. not all documents have a value

git_commits:

1. **summary:** LUCENE-7590: add DocValuesStatsCollector
message: LUCENE-7590: add DocValuesStatsCollector

github_issues:

github_issues_comments:

github_pulls:

github_pulls_comments:

github_pulls_reviews:

jira_issues:

1. **summary:** Add DocValues statistics helpers
description: I think it can be useful to have DocValues statistics helpers, that can allow users to query for the min/max/avg etc. stats of a DV field. In this issue I'd like to cover numeric DV, but there's no reason not to add it to other DV types too.
label: code-design
2. **summary:** Add DocValues statistics helpers
description: I think it can be useful to have DocValues statistics helpers, that can allow users to query for the min/max/avg etc. stats of a DV field. In this issue I'd like to cover numeric DV, but there's no reason not to add it to other DV types too.
3. **summary:** Add DocValues statistics helpers
description: I think it can be useful to have DocValues statistics helpers, that can allow users to query for the min/max/avg etc. stats of a DV field. In this issue I'd like to cover numeric DV, but there's no reason not to add it to other DV types too.
4. **summary:** Add DocValues statistics helpers
description: I think it can be useful to have DocValues statistics helpers, that can allow users to query for the min/max/avg etc. stats of a DV field. In this issue I'd like to cover numeric DV, but there's no reason not to add it to other DV types too.
5. **summary:** Add DocValues statistics helpers
description: I think it can be useful to have DocValues statistics helpers, that can allow users to query for the min/max/avg etc. stats of a DV field. In this issue I'd like to cover numeric DV, but there's no reason not to add it to other DV types too.
6. **summary:** Add DocValues statistics helpers
description: I think it can be useful to have DocValues statistics helpers, that can allow users to query for the min/max/avg etc. stats of a DV field. In this issue I'd like to cover numeric DV, but there's no reason not to add it to other DV types too.
label: code-design
7. **summary:** Add DocValues statistics helpers
description: I think it can be useful to have DocValues statistics helpers, that can allow users to query for the min/max/avg etc. stats of a DV field. In this issue I'd like to cover numeric DV, but there's no reason not to add it to other DV types too.
label: documentation
8. **summary:** Add DocValues statistics helpers
description: I think it can be useful to have DocValues statistics helpers, that can allow users to query for the min/max/avg etc. stats of a DV field. In this issue I'd like to cover numeric DV, but there's no reason not to add it to other DV types too.
label: code-design
9. **summary:** Add DocValues statistics helpers
description: I think it can be useful to have DocValues statistics helpers, that can allow users to query for the min/max/avg etc. stats of a DV field. In this issue I'd like to cover numeric DV, but there's no reason not to add it to other DV types too.
10. **summary:** Add DocValues statistics helpers
description: I think it can be useful to have DocValues statistics helpers, that can allow users to query for the min/max/avg etc. stats of a DV field. In this issue I'd like to cover numeric DV, but there's no reason not to add it to other DV types too.
11. **summary:** Add DocValues statistics helpers
description: I think it can be useful to have DocValues statistics helpers, that can allow users to query for the min/max/avg etc. stats of a DV field. In this issue I'd like to cover numeric DV, but there's no reason not to add it to other DV types too.
12. **summary:** Add DocValues statistics helpers

- [illegible]

not to add it to other DV types too.

41. **summary:** Add DocValues statistics helpers

description: I think it can be useful to have DocValues statistics helpers, that can allow users to query for the min/max/avg etc. stats of a DV field. In this issue I'd like to cover numeric DV, but there's no reason not to add it to other DV types too.

label: code-design

42. **summary:** Add DocValues statistics helpers

description: I think it can be useful to have DocValues statistics helpers, that can allow users to query for the min/max/avg etc. stats of a DV field. In this issue I'd like to cover numeric DV, but there's no reason not to add it to other DV types too.

43. **summary:** Add DocValues statistics helpers

description: I think it can be useful to have DocValues statistics helpers, that can allow users to query for the min/max/avg etc. stats of a DV field. In this issue I'd like to cover numeric DV, but there's no reason not to add it to other DV types too.

44. **summary:** Add DocValues statistics helpers

description: I think it can be useful to have DocValues statistics helpers, that can allow users to query for the min/max/avg etc. stats of a DV field. In this issue I'd like to cover numeric DV, but there's no reason not to add it to other DV types too.

jira_issues_comments:

1. **body:** First patch adds numeric statistics. I'd appreciate comments about it before I add support for sorted-numeric (including, whether we should!). Note that I chose to take either a field or `ValueSource`. The latter gives some flexibility by allowing users to pass an arbitrary VS over e.g. an `Expression` over a numeric DV field. This, as far as I could tell, does not apply to `SortedNumericDV`, or at least I couldn't find an existing `ValueSource` implementation (like `LongFieldSource`) for `SortedNumericDV`. If this approach looks good, I'd like to refactor the class so that it's easy to share/reuse code between Long and Double NDV fields.

label: code-design

2. Thanks [~shaie]. Maybe instead of a new stepping-stone class that the user must invoke, `DocsAndContexts`, you could just define a functional interface: `LeafReaderContext ctx` -> `DocIdSetIterator` And then maybe instead of a new class that computes stats in its ctor, `NumericLongDocValuesStats`, you could offer a static method instead, taking a top-level reader and the above function, that computes the stats and returns a results class holding the min/max/mean/etc.? Seems like that might be a simpler way to expose the functionality...
 3. I would not even define an own functional interface, just use `java.util.function.Function<LeafReaderContext,DocIdSetIterator>`: <https://docs.oracle.com/javase/8/docs/api/java/util/function/Function.html> This only works if the function does not throw checked exceptions. Otherwise a new functional interface is needed.
 4. Thanks [~mikemccand] and [~thetaphi], I changed to a static class and removed `DocsAndContexts` in favor of a new `Function<LeafReaderContext,DocIdSetIterator>`. Maybe `BitsDocIdSetIterator` can go in separately (i.e. a separate issue)? As I think it's a useful utility to have anyway.
 5. Thanks [~shaie] the patch looks great!
 6. **body:** bq. Maybe `BitsDocIdSetIterator` can go in separately (i.e. a separate issue)? It took us a lot of efforts to remove slow iterators so I'd like to not add them back. Let's implement the computation of these stats by writing a Collector and use a `MatchAllDocsQuery`? Why is `missing` undefined when `count` is zero?
- label:** code-design
7. **body:** bq. Let's implement the computation of these stats by writing a Collector and use a `MatchAllDocsQuery`? At first I thought this is an overkill, but a `Collector` will allow computing them for documents that match another query. I will explore that option. bq. Why is `missing` undefined when count is zero? I thought that if you have no documents in the index at all, then `missing` is undefined, but now that you ask the question, I guess in that case it's fine if it's `{0}`, like `count`. I'll change the docs.
- label:** documentation
8. **body:** Patch implements a `DocValuesStatsCollector`. Note some key design decisions: A `DocValuesStats` is responsible for providing the specific `DocValuesIterator` for a `LeafReaderContext`. It then accumulates the value, computes missing and other statistics. It computes `missing` and `count`, leaving `min` and `max` to the actual implementation. Also, this stats

does not define a `{mean}`, as at least for now I'm not sure how the mean value of a `{SortedSetDocValues}` is defined. An abstract `{NumericDocValuesStats}` implementation for single-numeric DV fields, which also adds a `{mean}` statistic, with two concrete implementations: `{LongNumericDocValuesStats}` and `{DoubleNumericDocValuesStats}`. This hierarchy should allow us to add further statistics for `{SortedSet}` and `{SortedNumeric}` DV fields. I did not implement them yet, as I'm not sure about some of the statistics (e.g. should the `{mean}` stat of a `{SortedNumeric}` be the mean across all values, or the minimum per document or ...). Let's discuss that separately. Also, note that I had to make `{DocValuesIterator}` public in order to declare it in this collector. If you're OK with the design and implementation, I want to separate `{DovValuesStats}` to its own file, for clarity. I did not do it yet though.

label: code-design

9. Added tests for `{DoubleNumericDocValuesStats}`. Now that I review the class names, how do you feel about removing `{Numeric}` from the concrete classes, so they're called `{Long/DoubleDocValuesStats}`?
10. Instead of using a `NOOP_COLLECTOR`, you could throw a `CollectionTerminatedException`, which will skip the segment entirely. By the way, in such cases I think we should still increase the missing count? Can we avoid making `{DocValuesIterator}` public?
11. bq. how do you feel about removing `Numeric` from the concrete classes, so they're called `Long/DoubleDocValuesStats`? Fine with me.
12. **body:** bq. Instead of using a `NOOP_COLLECTOR`, you could throw a `CollectionTerminatedException` OK, good idea. bq. By the way, in such cases I think we should still increase the missing count? I am not sure? I mean, `{missing}` represents all the documents that matched the query and did not have a value for that DV field. But when `{getLeafCollector}` is called, we don't know yet that any documents will be matched by the query at all (I think?) and therefore updating missing might be confusing? I.e., I'd expect that if anyone chained `{TotalHitsCollector}` with `{DocValuesStatsCollector}`, then `{totalHits = stats.count() + stats.missing()}`? I am open to discuss it, just not sure I always want to update missing with `{context.reader().numDocs()}` ... bq. Can we avoid making `DocValuesIterator` public? I did not find a way, since it's part of `{DocValuesStats.init()}` API and I think users should be able to provide their own `{Stats}` impl, e.g. if they want to compute something on a `{BinaryDocValues}` field? Here too, I'd love to get more ideas though. I tried to avoid implementing N collectors, one for each DV type, where they share a large portion of the code. But if you have strong opinions about making `{DVI}` public, maybe that's what we should do ...

label: code-design

13. **body:** [~jpountz] I accept your proposal about missing, only in case a reader does not have the requested DV field, the collector returns a `{LeafCollector}` which updates `{missing}` for every hit document. I also renamed the classes as proposed earlier, as well extracted `{DocValuesStats}` and friends to its own class. I still didn't address changing `{DocValuesIterator}` to public. BTW, I noticed that `{SimpleTextDocValuesReader}` defines a private class named `{DocValuesIterator}` with exactly the same signature, I assume because the other one is package-private. So I feel that changing `{DVI}` to public is beneficial beyond the scope of this issue alone. What do you think?

label: code-design

14. **body:** Thanks, the missing change and the rename look good to me. Regarding `DocValuesIterator`, I think in the simple text case, it is really an impl detail, and the one in core is only really used to share the declaration of `{advanceExact}`. My feeling is that making it public is only useful here because the way stats are computed is very abstracted with generics and inheritance, so I am not convinced this use-case requires that we make `DocValuesIterator` public. I would rather either avoid abstracting so much or define a functional interface in the stats package that only defines an `advanceExact` method and use method references to be able to share computation across the various doc value types. Maybe we can also make the collector pkg-private, it does not seem to need to be public, does it?

label: code-design

15. Patch changes `{DocValuesIterator}` package-private again and adds an API to `{DocValuesStats}` to help in determining whether a document has or does not have a value for the field. The Collector needs to be public because you're supposed to initialize it and run a search with it.
16. +1
17. Commit `ad7152ad4739a47aa2b45405ba1682b3dda18923` in lucene-solr's branch `refs/heads/master` from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=ad7152a>] LUCENE-7590: add `DocValuesStatsCollector`
18. Commit `43f4f7a279553913aadfdd684d9cdcff0a5f4220` in lucene-solr's branch `refs/heads/branch_6x` from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=43f4f7a>] LUCENE-7590: add

DocValuesStatsCollector

19. Commit e09ef681e4d36adb8987ca0cda6bcb3221830102 in lucene-solr's branch refs/heads/branch_6x from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=e09ef68>] Revert "LUCENE-7590: add DocValuesStatsCollector" This reverts commit 43f4f7a279553913aadfdd684d9cdcff0a5f4220.
20. Commit 7269c484a4a0dd147a445a4b676144592f0aa60f in lucene-solr's branch refs/heads/branch_6x from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=7269c48>] LUCENE-7590: add DocValuesStatsCollector
21. **body:** Commit 85582dabe4372085e1af5d01ebbfcd0303b9f12 in lucene-solr's branch refs/heads/master from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=85582da>] LUCENE-7590: fix typo in method parameter
label: documentation
22. There are now few tasks left: * Add more statistics, such as `{{sum}}` and `{{stdev}}` (for numeric fields). Should we care about overflow, or only document it? * We can also compute more stats like what Solr gives in [StatsComponent|<https://cwiki.apache.org/confluence/display/solr/The+Stats+Component#TheStatsComponent-StatisticsSupported>]. What do you think? * Add stats for `{{SortedDocValues}}`. This should be fairly straightforward by comparing the `{{BytesRef}}` of all matching documents. But I don't think we should have a `{{mean}}` stat for it? Likewise for `{{SortedSetDocValues}}`. * What should we do with `{{SortedNumericDocValues}}`? `{{min}}` and `{{max}}` are well defined, but what about `{{mean}}`? Should it be across all values? I intend to close this issue and handle the rest in follow-on issues, unless you think otherwise. Also, would appreciate your feedback on the above points.
23. Patch adds `{{sum}}`, `{{stdev}}` and `{{variance}}` stats to `{{NumericDocValuesStats}}`. I also added a CHANGES entry which I forgot to in the previous commit.
24. Commit 295cab7216ca76debaf4d354409741058a8641a1 in lucene-solr's branch refs/heads/master from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=295cab7>] LUCENE-7590: add sum, variance and stdev stats to NumericDVStats
25. Commit 2a0814fc34b76d8031938d09e11bedc7f604f543 in lucene-solr's branch refs/heads/branch_6x from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=2a0814f>] LUCENE-7590: add sum, variance and stdev stats to NumericDVStats
26. Patch adds DVStats for `{{SortedNumericDocValuesField}}`.
27. Commit 944b8e07f557b9320895998fe33d71cae5199eee in lucene-solr's branch refs/heads/master from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=944b8e0>] LUCENE-7590: add DocValuesStats for SortedNumeric DV fields
28. Commit 63a5cd00173f7e89a478981429d4d5cd38f3cf1d in lucene-solr's branch refs/heads/branch_6x from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=63a5cd0>] LUCENE-7590: add DocValuesStats for SortedNumeric DV fields
29. Patch adds `{{SortedDocValuesStats}}` and `{{SortedSetDocValuesStats}}` for sorted and sorted-set DV fields. With this patch, I think the issue is ready to be closed. I am not sure that we need a DVStats for a BinaryDVField at this point, but if demand arises, it should be easy to add.
30. Commit 23206caabd09310cb23a2b5302ce41af62b5c9cc in lucene-solr's branch refs/heads/master from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=23206ca>] LUCENE-7590: add Sorted(Set)DocValuesStats
31. Commit 47bb32c3bb77a2dfaaf9d1db50e244599cf053a6 in lucene-solr's branch refs/heads/branch_6x from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=47bb32c>] LUCENE-7590: add Sorted(Set)DocValuesStats
32. Commit 73b6a29f2d89e2f1ce86b57ad0acac7d157f7e21 in lucene-solr's branch refs/heads/branch_6x from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=73b6a29>] LUCENE-7590: move docsWithField to DocValuesStats
33. Committed to master and 6x. This is now complete.
34. Commit 321c6f090f04463a8798d090e5426efeabdbdc418 in lucene-solr's branch refs/heads/master from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=321c6f0>] LUCENE-7590: make (Sorted)NumericDocValuesStats public
35. Commit f075a673c9629e92c1e9dd1e104a4e602d6fe610 in lucene-solr's branch refs/heads/branch_6x from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=f075a67>] LUCENE-7590: make (Sorted)NumericDocValuesStats public
36. Commit 944b8e07f557b9320895998fe33d71cae5199eee in lucene-solr's branch refs/heads/feature/metrics from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=944b8e0>] LUCENE-7590: add DocValuesStats for SortedNumeric DV fields

37. Commit 23206caabd09310cb23a2b5302ce41af62b5c9cc in lucene-solr's branch refs/heads/feature/metrics from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=23206ca>] LUCENE-7590: add Sorted(Set)DocValuesStats
38. Commit 321c6f090f04463a8798d090e5426efeabbdc418 in lucene-solr's branch refs/heads/feature/metrics from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=321c6f0>] LUCENE-7590: make (Sorted)NumericDocValuesStats public
39. Commit c083e81e6015d8d52ccd74ad1e966862936fb926 in lucene-solr's branch refs/heads/branch_6x from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=c083e81>] LUCENE-7590: fix test edge case In case all indexed documents were deleted, the test failed to correctly assert the number of expected missing documents.
40. Commit 4d81eee8a141c68b17c2f75cf6534fb352d94473 in lucene-solr's branch refs/heads/master from [~shaie] [<https://git-wip-us.apache.org/repos/asf?p=lucene-solr.git;h=4d81eee>] LUCENE-7590: fix test edge case In case all indexed documents were deleted, the test failed to correctly assert the number of expected missing documents.
41. **body:** [~shaie] I'm a little confused about the description defined here. {code} public LongDocValuesStats(String description) { super(description, Long.MAX_VALUE, Long.MIN_VALUE); } {code} it was in turn passed to NumericDocValueStats as the name of DV field. Why not use the {{field}} as the name of parameter in LongDocValuesStats?
label: code-design
42. [~shia] where do you see that? I checked master and there's no {{description}} in the file at all. Here's the code: {code} public LongDocValuesStats(String field) { super(field, Long.MAX_VALUE, Long.MIN_VALUE); } {code}