Articles → 591. Tag Validator ▼

(1) (1) (ii)

## 591. Tag Validator

June 3, 2017 | 10K views

\*\*\* Average Rating: 4.27 (11 votes)

Given a string representing a code snippet, you need to implement a tag validator to parse the code and return whether it is valid. A code snippet is valid if all the following rules hold:

- 1. The code must be wrapped in a valid closed tag. Otherwise, the code is invalid. 2. A closed tag (not necessarily valid) has exactly the following format:
- <TAG\_NAME>TAG\_CONTENT</TAG\_NAME> . Among them, <TAG\_NAME> is the start tag, and </TAG NAME> is the end tag. The TAG\_NAME in start and end tags should be the same. A closed tag is valid if and only if the TAG\_NAME and TAG\_CONTENT are valid. 3. A valid TAG NAME only contain upper-case letters, and has length in range [1,9]. Otherwise, the TAG NAME is invalid.
- 4. A valid TAG CONTENT may contain other valid closed tags, cdata and any characters (see note1) **EXCEPT** unmatched < , unmatched start and end tag, and unmatched or closed tags with invalid TAG\_NAME. Otherwise, the TAG\_CONTENT is invalid.
- 5. A start tag is unmatched if no end tag exists with the same TAG\_NAME, and vice versa. However, you also need to consider the issue of unbalanced when tags are nested. 6. A < is unmatched if you cannot find a subsequent > . And when you find a < or </ , all the
- subsequent characters until the next > should be parsed as TAG\_NAME (not necessarily valid). 7. The cdata has the following format: <![CDATA[CDATA CONTENT]]>. The range of CDATA CONTENT is
- defined as the characters between <![CDATA[ and the first subsequent ]]>. 8. CDATA\_CONTENT may contain any characters. The function of cdata is to forbid the validator to parse CDATA\_CONTENT, so even it has some characters that can be parsed as tag (no matter valid or invalid),
- you should treat it as regular characters.
- Valid Code Examples: Input: "<DIV>This is the first line <![CDATA[<div>]]></DIV>"

## Output: True

```
Explanation:
The code is wrapped in a closed tag : <DIV> and </DIV>.
The TAG_NAME is valid, the TAG_CONTENT consists of some characters and cdata.
Although CDATA_CONTENT has unmatched start tag with invalid TAG_NAME, it should be con
So TAG_CONTENT is valid, and then the code is valid. Thus return true.
Input: "<DIV>>> ![cdata[]] <![CDATA[<div>]>]]>>]</DIV>"
Output: True
Explanation:
We first separate the code into : start_tag|tag_content|end_tag.
start tag -> "<DIV>"
end_tag -> "</DIV>"
tag_content could also be separated into : text1|cdata|text2.
text1 -> ">> ![cdata[]] "
cdata -> "<![CDATA[<div>]>]]>", where the CDATA_CONTENT is "<div>]>"
text2 -> "]]>>]"
The reason why start_tag is NOT "<DIV>>>" is because of the rule 6.
The reason why cdata is NOT "<![CDATA[<div>]>]]>" is because of the rule 7.
```

## Output: False

Output: False

Invalid Code Examples:

Input: "<A> <B> </A> </B>"

```
Input: "<DIV> div tag is not closed <DIV>"
 Input: "<DIV> unmatched < </DIV>"
 Output: False
 Input: "<DIV> closed tags with invalid tag name <b>123</b> </DIV>"
 Output: False
 Input: "<DIV> unmatched tags with invalid tag name </1234567890> and <CDATA[[]]> 
 Output: False
 Input: "<DIV> unmatched start tag <B> and unmatched end tag </C> </DIV>"
 Output: False
Note:

    For simplicity, you could assume the input code (including the any characters mentioned above) only

    contain letters, digits, '<','>','/','!','[',']' and ' '.
```

Explanation: Unbalanced. If "<A>" is closed, then "<B>" must be unmatched, and vice ve

Solution

Approach 1: Stack

checking the following properties.

 The code should be wrapped in valid closed tag. The TAG\_NAME should be valid.

Summarizing the given problem, we can say that we need to determine whether a tag is valid or not, by

5. All the tags should be closed, i.e. each start-tag should have a corresponding end-tag and vice-versa and the order of the tags should be correct as well.

4. The cdata should be valid.

3. The TAG CONTENT should be valid.

- In order to check the validity of all these, firstly, we need to identify which parts of the given code string act as which part from the above mentioned categories. To understand how it's done, we'll go through the
- implementation and the reasoning behind it step by step. We iterate over the given code string. Whenever a < is encountered(unless we are currently inside <!

the code string turns out to be invalid as well.

last end-tag has been encountered.

Java

9

10

11

12 13

14 15

16

17

18

19 20

21

22

1 public class Solution {

Stack < String > stack = new Stack < > ();

if (s.length() < 1 || s.length() > 9)

for (int i = 0; i < s.length(); i++) {

return false;

stack.pop();

return false;

contains\_tag = true;

public boolean isValidCdata(String s) {

{min,} The preceding item is matched min or more times.

[...] Matches any single character in brackets.

[^...] Matches any single character not in brackets.

of the inner tags) by using the regex expression below:

letters with length between 1 to 9 lying inside <...> or </..> 4: Matches cdata. Any characters enclosed within <![CDATA[...]]>

cdata characters only till finding the first ]]>, we need to use .\*?

only need to check the presence of inner closed tags.

Check this link for testing any regular expression on a sample text.

public boolean isValidTagName(String s, boolean ending) {

if (!stack.isEmpty() && stack.peek().equals(s))

Stack < String > stack = new Stack < > ();

stack.pop();

return false;

contains\_tag = true;

0 ∧ ∨ ☑ Share ← Reply

boolean contains\_tag = false;

if (ending) {

else

} else {

outermost start-tag.

Java

10

11

12

1 import java.util.regex.\*; 2 public class Solution {

stack.push(s);

public boolean isValidTagName(String s, boolean ending) {

if (!Character.isUpperCase(s.charAt(i)))

if (!stack.isEmpty() && stack.peek().equals(s))

boolean contains\_tag = false;

return false;

if (ending) {

} else {

return true;

The following animation depicts the process.

cdata as per the conditions given in the problem statement.

If the character immediately following this < is an !, the characters following this < can't be a part of a valid TAG NAME, since only upper-case letters(in case of a start tag) or / followed by upper-case letters(in the case of an end tag). Thus, the choice now narrows down to only cdata. Thus, we need to check if the

current bunch of characters following <! (including it) constitute a valid cdata. For doing this, firstly we find out the first matching 11> following the current <! to mark the ending of cdata. If no such matching 11> exists, the code string is considered as invalid. Apart from this, the <! should also be immediately followed

[CDATA[...]]> ), it indicates the beginning of either a TAG\_NAME (start tag or end tag) or the beginning of

by CDATA[ for the cdata to be valid. The characters lying inside the <![CDATA[ and ]]> do not have any constraints on them. If the character immediately following the < encountered isn't an !, this < can only mark the beginning of TAG NAME. Now, since a valid start tag can't contain anything except upper-case letters, if a / is found after < , the </ pair indicates the beginning of an end tag. Now, when a < refers to the beginning of a

TAG\_NAME (either start-tag or end-tag), we find out the first closing > following the < to find out the substring(say s), that constitutes the TAG\_NAME. This s should satisfy all the criterion to constitute a valid TAG\_NAME. Thus, for every such s, we check if it contains all upper-case letters and also check its length(It should be between 1 to 9). If any of the criteria isn't fulfilled, s doesn't constitue a valid TAG\_NAME. Hence,

Apart from checking the validity of the TAG\_NAME, we also need to ensure that the tags always exist in pairs. i.e. for every start-tag, a corresponding end-tag should always exist. Further, we can note that in case of multiple TAG\_NAME 's, the TAG\_NAME whose start-tag comes later than the other ones, should have its endtag appearing before the end-tags of those other TAG\_NAME 's. i.e. the tag which starts later should end first. From this, we get the intuition that we can make use of a stack to check the existence of matching start and end-tags. Thus, whenever we find out a valid start-tag, as mentioned above, we push its TAG\_NAME string onto a stack. Now, whenever an end-tag is found, we compare its TAG NAME with the TAG NAME at the

top the stack and remove this element from the stack. If the two don't match, this implies that either the current end-tag has no corresponding start-tag or there is a problem with the ordering of the tags. The two need to match for the tag-pair to be valid, since there can't exist an end-tag without a corresponding starttag and vice-versa. Thus, if a match isn't found, we can conclude that the given code string is invalid.

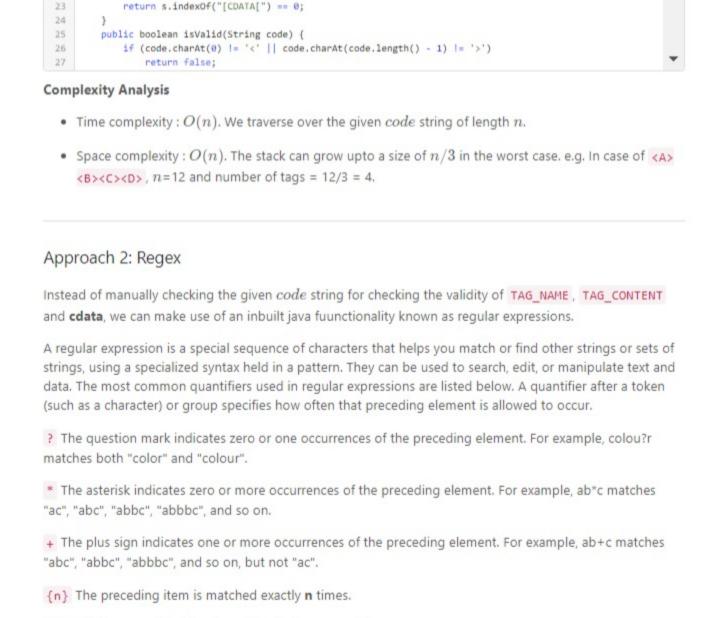
Now, after the complete code string has been traversed, the stack should be empty if all the start-tags have their corresponding end-tags as well. If the stack isn't empty, this implies that some start-tag doesn't have the corresponding end-tag, violating the closed-tag's validity condition. Further, we also need to ensure that the given code is completely enclosed within closed tags. For this, we need to ensure that the first cdata found is also inside the closed tags. Thus, when we find a possibility of the

presence of cdata, we proceed further only if we've already found a start tag, indicated by a non-empty stack. Further, to ensure that no data lies after the last end-tag, we need to ensure that the stack doesn't become empty before we reach the end of the given code string, since an empty stack indicates that the

< X Y Z > o n e < B > T w 0 < ! [ C D A T A [ t h < ! e ] ] > 4 < / B > < / X Y Z >

1/21

**Сору** 



 $<([A-Z]{1,9})>([^<]*((<\/?[A-Z]{1,9}>)|(<!\[CDATA\[(.*?)]]>))?[^<]*)*<\/\1>$ The image below shows the portion of the string that each part of the expression helps to match:

1: Matches the outermost start-tag. All upper-case alphabets with length between 1 to 9 lying inside <..>

3 : Matches all TAG\_DATA(both start-tag and end-tag) lying inside the outer closed tags. All upper-case

5: Matches the outermost end-tag. All upper-case alphabets with length between 1 to 9 lying inside </..> Uses back-reference with the first parentheses. Thus, the outermost end-tag should be the same as the

6: By default, .\* is a greedy search and matches as much characters as possible. To limit the match of

2: Matches the TAG\_CONTENT(except cdata). All characters except < occurring 0 or more times.

Thus, by making use of regex, we can directly check the validity of the code string directly(except the nesting

{min, max} The preceding item is matched at least min times, but not more than max times.

() Parentheses are used to define the scope and precedence of the operators (among other uses). For example, gray|grey and gr(a|e)y are equivalent patterns which both describe the set of "gray" or "grey".

A vertical bar separates alternatives. For example, gray|grey can match "gray" or "grey".

But, if we make use of back-referencing as mentioned above, the matching process takes a very large amount of CPU time. Thus, we use the regex only to check the validity of the TAG\_CONTENT, TAG\_NAME and the **cdata**. We check the presence of the outermost closed tags by making use of a stack as done in the last approach. The rest of the process remains the same as in the last approach, except that we need not manually check

the validity of TAG\_CONTENT, TAG\_NAME and the cdata, since it is already done by the regex expression. We

Copy Copy

13 stack.push(s); 14 15 16 17 public boolean isValid(String code) { 18 if (!Pattern.matches(regex, code))

19 return false; for (int i = 0; i < code.length(); i++) { boolean ending = false; 22 23 if (stack.isEmpty() && contains\_tag) return false: 24 if (code.charAt(i) == '<') { 26 if (code.charAt(i + 1) == '!') { 27 i = code.indexOf("]]>", i + 1); Complexity Analysis Time complexity: Regular Expressions are/can be implemented in the form of Finite State Machines. Thus, the time complexity is dependent on the internal representation. In case of any suggestions, please comment below. • Space complexity : O(n). The stack can grow upto a size of n/3 in the worst case. e.g. In case of  $\langle A \rangle$  $\langle B \rangle \langle C \rangle \langle D \rangle$ , n=12 and number of tags = 12/3=4. Rate this article: \* \* \* \* \* 3 Previous Next **1** Comments: 4 Sort By ▼

