

Why doesn't `b` match word boundaries correctly?



Raymond Chen

A colleague of mine was having trouble getting the `\b` metacharacter in a regular expression to work. Of course, when somebody asks a question like that, you first have to establish what their definition of “work” is. Fortunately, he provided some examples:

<code>Regex.IsMatch("foo", @"\b" + "foo" + @"\b")</code>	true
<code>Regex.IsMatch("%1" , @"\b" + "%1" + @"\b")</code>	false
<code>Regex.IsMatch("%1" , @"\b" + @"\%1" + @"\b")</code>	false
<code>Regex.IsMatch("%1" , @"\b" + @"\%1" + @"\b")</code>	false
<code>Regex.IsMatch("%1" , @".")</code>	true
<code>Regex.IsMatch("%1" , @"%1")</code>	true

“The last two entries are just sanity checks to make sure I didn’t make some stupid mistake like passing the parameters in the wrong order. I want to search for a string that contains `%1` with word boundaries on either side, something I would normally use `\b` for. Is there something special about the `%` character? Notice that the match succeeds when I look for the word `foo`.” Everything is working as it should. Recall that the `\b` metacharacter matches when there is a `\w` on one side and a `\W` on the other, where the beginning and end of the string are treated as if they were `\w`. The string `%1` therefore breaks down as

virtual <code>\w</code>	beginning of string
<code>\w</code>	<code>%</code> is not an alphanumeric or <code>_</code>
<code>\w</code>	<code>1</code> is a digit
virtual <code>\W</code>	end of string

The only points where `\b` would match are immediately before and after the `1`, since those are the transition points between `\w` and `\W` and vice versa. In particular, the location immediately before the percent sign does not match since it is surrounded by `\W` on both sides.

My colleague responded, “D’oh! I keep forgetting that `%` won’t act like a `\w` just because I want it to.”

Raymond Chen

Follow

