
**Information technology — Guidance
for the use of database language
SQL —**

**Part 1:
XQuery regular expressions**

*Technologies de l'information — Recommandations pour l'utilisation
du langage de base de données SQL —*

Partie 1: Expressions régulières de XQuery en SQL

IECNORM.COM : Click to view the full PDF of ISO/IEC 19075-1:2021



IECNORM.COM : Click to view the full PDF of ISO/IEC 19075-1:2021



COPYRIGHT PROTECTED DOCUMENT

© ISO/IEC 2021

All rights reserved. Unless otherwise specified, or required in the context of its implementation, no part of this publication may be reproduced or utilized otherwise in any form or by any means, electronic or mechanical, including photocopying, or posting on the internet or an intranet, without prior written permission. Permission can be requested from either ISO at the address below or ISO's member body in the country of the requester.

ISO copyright office
CP 401 • Ch. de Blandonnet 8
CH-1214 Vernier; Geneva
Phone: +41 22 749 01 11
Email: copyright@iso.org
Website: www.iso.org

Published in Switzerland

| Contents | Page |
|--|-------------|
| Foreword..... | v |
| Introduction..... | vii |
| 1 Scope..... | 1 |
| 2 Normative references..... | 2 |
| 3 Terms and definitions..... | 3 |
| 4 XQuery regular expressions..... | 4 |
| 4.1 Context of XQuery regular expressions..... | 4 |
| 4.2 Introduction to XQuery regular expressions..... | 4 |
| 4.3 Matching a specific character..... | 4 |
| 4.4 Metacharacters and escape sequences..... | 5 |
| 4.5 Dot..... | 6 |
| 4.6 Anchors..... | 7 |
| 4.7 Line terminators..... | 7 |
| 4.8 Bracket expressions..... | 8 |
| 4.8.1 Introduction to bracket expressions..... | 8 |
| 4.8.2 Listing characters..... | 8 |
| 4.8.3 Matching a range..... | 9 |
| 4.8.4 Negation..... | 9 |
| 4.8.5 Character class subtraction..... | 9 |
| 4.9 Alternation..... | 9 |
| 4.10 Quantifiers..... | 10 |
| 4.11 Locating a match..... | 11 |
| 4.12 Capture and back-reference..... | 12 |
| 4.13 Precedence..... | 13 |
| 4.14 Modes..... | 13 |
| 5 Operators using regular expressions..... | 15 |
| 5.1 Introduction to operators using regular expressions..... | 15 |
| 5.2 LIKE_REGEX..... | 15 |
| 5.3 OCCURRENCES_REGEX..... | 16 |
| 5.4 POSITION_REGEX..... | 17 |
| 5.5 SUBSTRING_REGEX..... | 19 |
| 5.6 TRANSLATE_REGEX..... | 20 |
| Bibliography..... | 23 |
| Index..... | 24 |

Tables

| Table | Page |
|-------------------------|------|
| 1 Match priorities..... | 11 |

IECNORM.COM : Click to view the full PDF of ISO/IEC 19075-1:2021

Foreword

ISO (the International Organization for Standardization) and IEC (the International Electrotechnical Commission) form the specialized system for worldwide standardization. National bodies that are members of ISO or IEC participate in the development of International Standards through technical committees established by the respective organization to deal with particular fields of technical activity. ISO and IEC technical committees collaborate in fields of mutual interest. Other international organizations, governmental and non-governmental, in liaison with ISO and IEC, also take part in the work.

The procedures used to develop this document and those intended for its further maintenance are described in the ISO/IEC Directives, Part 1. In particular, the different approval criteria needed for the different types of document should be noted. This document was drafted in accordance with the editorial rules of the ISO/IEC Directives, Part 2 (see www.iso.org/directives or www.iec.ch/members_experts/refdocs).

Attention is drawn to the possibility that some of the elements of this document may be the subject of patent rights. ISO and IEC shall not be held responsible for identifying any or all such patent rights. Details of any patent rights identified during the development of the document will be in the Introduction and/or on the ISO list of patent declarations received (see www.iso.org/patents) or the IEC list of patent declarations received (see patents.iec.ch).

Any trade name used in this document is information given for the convenience of users and does not constitute an endorsement.

For an explanation of the voluntary nature of standards, the meaning of ISO specific terms and expressions related to conformity assessment, as well as information about ISO's adherence to the World Trade Organization (WTO) principles in the Technical Barriers to Trade (TBT) see www.iso.org/iso/foreword.html. In the IEC, see www.iec.ch/understanding-standards.

This document was prepared by Joint Technical Committee ISO/IEC JTC 1, *Information technology*, Subcommittee SC 32, *Data management and interchange*.

This first edition of ISO/IEC 19075-1 cancels and replaces ISO/IEC TR 19075-1:2011.

This document is intended to be used in conjunction with the following editions of the parts of the ISO/IEC 9075 series:

- ISO/IEC 9075-1, sixth edition or later;
- ISO/IEC 9075-2, sixth edition or later;
- ISO/IEC 9075-3, sixth edition or later;
- ISO/IEC 9075-4, seventh edition or later;
- ISO/IEC 9075-9, fifth edition or later;
- ISO/IEC 9075-10, fifth edition or later;
- ISO/IEC 9075-11, fifth edition or later;
- ISO/IEC 9075-13, fifth edition or later;

ISO/IEC 19075-1:2021(E)

- ISO/IEC 9075-14, sixth edition or later;
- ISO/IEC 9075-15, second edition or later;
- ISO/IEC 9075-16, first edition or later.

A list of all parts in the ISO/IEC 19075 series can be found on the ISO and IEC websites.

Any feedback or questions on this document should be directed to the user's national standards body. A complete listing of these bodies can be found at www.iso.org/members.html and www.iec.ch/national-committees.

IECNORM.COM : Click to view the full PDF of ISO/IEC 19075-1:2021

Introduction

The organization of this document is as follows:

- 1) **Clause 1, “Scope”**, specifies the scope of this document.
- 2) **Clause 2, “Normative references”**, identifies additional standards that, through reference in this document, constitute provisions of this document.
- 3) **Clause 3, “Terms and definitions”**, defines the terms and definitions used in this document.
- 4) **Clause 4, “XQuery regular expressions”**, explains how XQuery regular expressions are formed.
- 5) **Clause 5, “Operators using regular expressions”**, explains how the SQL operators use regular expressions.

IECNORM.COM : Click to view the full PDF of ISO/IEC 19075-1:2021

[IECNORM.COM](https://www.iecnorm.com) : Click to view the full PDF of ISO/IEC 19075-1:2021

Information technology — Guidance for the use of database language SQL —

Part 1:

XQuery regular expressions

1 Scope

This document describes the regular expression support in SQL (ISO/IEC 9075-2) adopted from the regular expression syntax of XQuery and XPath Functions and Operators 3.1, which is derived from Perl. This document discusses five operators using this regular expression syntax:

- LIKE_REGEX predicate, to determine the existence of a match to a regular expression.
- OCCURRENCES_REGEX numeric function, to determine the number of matches to a regular expression.
- POSITION_REGEX function, to determine the position of a match.
- SUBSTRING_REGEX function, to extract a substring matching a regular expression.
- TRANSLATE_REGEX function, to perform replacements using a regular expression.

IECNORM.COM : Click to view the full PDF of ISO/IEC 19075-1:2021

2 Normative references

The following documents are referred to in the text in such a way that some or all of their content constitutes requirements of this document. For dated references, only the edition cited applies. For undated references, the latest edition of the referenced document (including any amendments) applies.

ISO/IEC 9075-1, *Information technology — Database languages — SQL — Part 1: Framework (SQL/Framework)*

ISO/IEC 9075-2, *Information technology — Database languages — SQL — Part 2: Foundation (SQL/Foundation)*

ISO/IEC 9075-14, *Information technology — Database languages — SQL — Part 14: XML-Related Specifications (SQL/XML)*

Bray, Tim et al.. *Extensible Markup Language (XML) Version 1.0, W3C Recommendation* [online]. Fifth Edition. Cambridge, Massachusetts, USA: W3C, 26 November 2008. Available at <http://www.w3.org/TR/xml>

Bray, Tim et al. *Extensible Markup Language (XML) Version 1.1, W3C Recommendation* [online]. Second Edition. Cambridge, Massachusetts, USA: W3C, 16 August 2006. Available at <http://www.w3.org/TR/xml11>

Biron, Paul V.; Malhotra, Ashok. *XML Schema Part 2: Datatypes, W3C Recommendation* [online]. Second Edition. Cambridge, Massachusetts, USA: W3C, 28 October 2004. Available at <http://www.w3.org/TR/xmlschema-2/>

Malhotra, Ashok et al.. *XQuery and XPath Functions and Operators 3.1, W3C Recommendation* [online]. Cambridge, Massachusetts, USA: W3C, 21 March 2017. Available at <http://www.w3.org/TR/xpath-functions/>

The Unicode Consortium. *Unicode Regular Expressions* [online]. 21. Mountain View, California, USA: The Unicode Consortium, 2020-06-17. Available at <http://www.unicode.org/reports/tr18/tr18-21.html>

3 Terms and definitions

For the purposes of this document, the terms and definitions given in ISO/IEC 9075-1 apply.

ISO and IEC maintain terminological databases for use in standardization at the following addresses:

- IEC Electropedia: available at <http://www.electropedia.org/>
- ISO Online browsing platform: available at <http://www.iso.org/obp>

IECNORM.COM : Click to view the full PDF of ISO/IEC 19075-1:2021

4 XQuery regular expressions

4.1 Context of XQuery regular expressions

The requirements for the material discussed in this document shall be as specified in ISO/IEC 9075-2, ISO/IEC 9075-14, XML 1.0, XML 1.1, XML Schema Part 2: Datatypes, XQuery and XPath Functions and Operators 3.1, and Unicode Technical Standard #18.

4.2 Introduction to XQuery regular expressions

This document explains the manner in which XQuery regular expressions are used by database language SQL in ISO/IEC 9075-2 and in ISO/IEC 9075-14. Both ISO/IEC 9075-2 and ISO/IEC 9075-14 specify requirements for the material discussed in this document.

XQuery regular expression syntax is specified in XQuery and XPath Functions and Operators 3.1, section 5.6.1, “Regular expression syntax”. This paper references the XQuery specification, with two small modifications (required since character strings in an RDBMS are not necessarily normalized according to XML conventions). The following subsections provide an overview of this syntax.

The XQuery regular expression syntax is itself a modification of another regular expression syntax found in XML Schema Part 2: Datatypes.

This section presents an overview of the capabilities of XQuery regular expression syntax. In the process, this section will illustrate some of the SQL operators. The SQL operators themselves are presented in Clause 5, “Operators using regular expressions”.

The following discussion does not cover every aspect of XQuery regular expressions; for this, XQuery and XPath Functions and Operators 3.1 is the reference (though hardly a tutorial; a variety of popular works contain detailed treatments of regular expressions).

4.3 Matching a specific character

Perhaps the most elementary pattern matching requirement is the ability to match a single character or string. For most characters, this is done by simply writing the character in the regular expression. For example, suppose an application needs to know if a string *S* contains the letters “xyz”. This could be done with the following predicate:

```
S LIKE_REGEX 'xyz'
```

Note that the SQL LIKE predicate would require an exact match for “xyz”. However, the convention with regular expressions is that *S* need only contain a substring that is “xyz”. For example, all of the following values of *S* would yield *True* for the immediately preceding predicate:

```
xyz
abcxyz123
1 xyz 2 xyz 3 xyz
```

Note that in the last example, there are actually three occurrences of the regular expression “xyz” within the tested value. The user may wish to know the number of occurrences of a match. This can be done with OCCURRENCES_REGEX. For example:

```
OCCURRENCES_REGEX ('xyz' IN '1 xyz 2 xyz 3 xyz') = 3
```

The user might also wish to know the position of a specific match. This can be done using POSITION_REGEX. For example, to learn the starting character position of the second occurrence,

```
POSITION_REGEX ('xyz' IN '1 xyz 2 xyz 3 xyz' OCCURRENCE 2 ) = 9
```

It is also possible to ask for the character position of the first character after the match. For example:

```
POSITION_REGEX ( AFTER 'xyz' IN '1 xyz 2 xyz 3 xyz' OCCURRENCE 2 ) = 12
```

If AFTER is used and the last character of the subject string is consumed, then the result is the length of the string plus 1 (one):

```
POSITION_REGEX ( AFTER 'xyz' IN 'xyz' ) = 4
```

4.4 Metacharacters and escape sequences

As mentioned, most characters can be matched by simply writing the character in the regular expression. However, certain characters are reserved as *metacharacters*. The complete list of metacharacters is:

```
. \ ? * + { } ( ) | [ ] ^ $
```

The use of each of these metacharacters is explained in subsequent Subclauses. To match a metacharacter, then an *escape sequence*, consisting of a backslash (“\”) followed by the metacharacter is used. For example, this expression tests whether a string contains a dollar sign:

```
S LIKE_REGEX '\$'
```

In particular, the escape sequence representing a backslash is two consecutive backslashes. There are various other defined escape sequences, matching either a single character, or any of a group of characters. The *single character escape sequences* are:

- \n newline (U+000A)
- \r return (U+000D)
- \t tab (U+0009)
- \- minus sign ('-')

The so-called *category escapes* are exemplified by “\p{L}” or “\p{Lu}”. A category escape begins with “\p{” followed by one uppercase letter, optionally a lowercase letter, and then the closing brace. In these example, “\p{L}” matches any letter (as defined by Unicode) and “\p{Lu}” matches any uppercase letter. Some interesting category escapes are:

- \p{L} Any letter.
- \p{Lu} Any uppercase letter.
- \p{Ll} Any lowercase letter.
- \p{Nd} Any decimal digit.
- \p{P} Any punctuation mark.

4.4 Metacharacters and escape sequences

`\p{Z}` Any separator (space, line, paragraph, etc.).

The complete list of category escapes is found in [XML Schema Part 2: Datatypes](#), section F.1.1, “Character class escapes”.

There are also *complementary category escapes*, which are exemplified by “`\P{L}`” or “`\P{Lu}`”. A complementary category escape matches any character that would not be matched by the corresponding category escape. The difference is that the (positive) character escape is written with a lowercase “p” whereas the complementary character escape is written with an uppercase “P”.

The so-called *block escapes* match any character in a block of Unicode, that is, a predefined consecutive range of code points. For example, “`\p{IsBasicLatin}`” matches the ASCII character set. There are also *complementary block escapes*, such as “`\P{IsBasicLatin}`”, which matches any single character that is not an ASCII character.

Finally, there are the following *multi-character escape sequences*:

- `\s` As defined by [XML Schema Part 2: Datatypes](#), this escape matches space (U+0020), tab (U+0009), newline (U+000A), or return (U+000D). Since character strings in an RDBMS have not undergone XML line termination normalization, the escape is broadened to include any character or two-character sequence that is recognized by [Unicode Technical Standard #18](#) as a line terminator. [Subclause 4.7, “Line terminators”](#), discusses this issue further.
- `\S` Any single character not matched by `\s`.
- `\i` Underscore (“_”), colon (“:”) or letter (this is a lot more than just the Latin letters; see [XML 1.0](#) appendix B, rule [84]).
- `\I` Any single character not matched by `\i`.
- `\c` Any single character matched by `NameChar`, as defined in [XML 1.0](#) section 2.3, rule [4a].
- `\C` Any single character not matched by `\c`.
- `\d` Any single digit
- `\D` Any single character not matched by `\d`.
- `\w` Any single Unicode character except those classified as “punctuation”, “separator”, or “other”.
- `\W` The complement of `\w`.

4.5 Dot

Dot (period, “.”) is a metacharacter that is used to match any single character (the same behavior as “_” in LIKE predicates), or any single character that is not a line terminator. The default is to match anything except a line terminator. The alternative, called *dot-all mode*, is specified using a flag that contains a lowercase “s”.

For example

```
S LIKE_REGEX 'a.b'
```

matches the following:

```
'xa0by'
```

but not the following:

```
'xa
by'
```

because the character between the “a” and the “b” is a line terminator. However, using dot-all mode like this:

```
S LIKE_REGEX 'a.b' FLAG 's'
```

would match both examples.

4.6 Anchors

As discussed in this Clause, regular expressions look for a match anywhere within a string, without needing to match the entire string. But how can an application require a match of the entire string? For this, the application uses *anchors*. The anchors are the metacharacters “^” for the start of a string (or line), and “\$” for the end of a string (or line). For example:

```
S LIKE_REGEX '^xyz$'
```

can only match a string that is precisely 'xyz'.

Anchors may be used separately to require a “begins with” or “end with” match. For example

```
S LIKE_REGEX '^xyz'
```

matches any string that begins with “xyz”, and

```
S LIKE_REGEX 'xyz$'
```

matches any string that ends with “xyz”.

Instead of matching the begin or end of the string, the anchors may be used to anchor a match to the begin or end of a line, by performing the match in *multi-line mode*. Multi-line mode is specified using a flag containing a lowercase “m”. For example:

```
S LIKE_REGEX '^xyz' FLAG 'm'
```

performs an anchored search in multi-line mode, matching any string containing a line that begins with “xyz”. The immediately preceding example would match the following string:

```
'line one
xyz
line three'
```

4.7 Line terminators

The metacharacters “.”, “^”, and “\$” and the multi-character escape sequences “\s” and “\S” are defined in terms of a “line terminator”. What counts as a line terminator? [XQuery and XPath Functions and Operators 3.1](#) only recognizes a line feed (U+000A) as a line terminator. This definition works well for XQuery, because XML normalizes the line terminators on various platforms to a line feed.

A closer look shows that XML has two definitions of line handling, in section 2.11, “End-of-line handling”, of [XML 1.0](#) and [XML 1.1](#). So which should be used for SQL?

4.7 Line terminators

A first step in answering this is to look at the definition of <XML query> in ISO/IEC 9075-14, which requires XML 1.0 as a basic level of support, and permits XML 1.1 support in the form of Feature X211, “XML 1.1 support”. So, character string can be normalized according to either XML 1.0 or XML 1.1 as an implementation-defined choice, or perhaps via a conformance feature.

However, some of the line terminators, even in XML 1.0, are two-character sequences. XML normalizes its input, which means that such two-character sequences are converted to a single character. This changes the relative position of every subsequent character, which would cause unexpected results for POSITION_REGEX.

Our solution is to look to [Unicode Technical Standard #18](#), a Unicode standard containing guidelines for regular expression processors. This provides a referenceable definition of line terminator that does not require normalizing the subject character string.

4.8 Bracket expressions

4.8.1 Introduction to bracket expressions

So far, this Clause has shown how to match a specific character or any character from certain predefined sets of characters. Using bracket expressions, applications can specify user-defined group of characters. (XML Schema and XQuery call these *character class expressions*, but the term *bracket expression* is in common use.)

A bracket expression is begun by a left bracket “[” and terminated by a right bracket “]”. Bracket expressions have a different list of special characters, namely

`^ [] \`

For clarity, these are called *special characters*, in contrast to the metacharacters listed earlier.

4.8.2 Listing characters

If a bracket expression does not contain any of the special characters, then the bracket expression matches any single character that is listed between the brackets. For example,

```
S LIKE_REGEX '[abc]'
```

matches any of the following:

```
'say'
'boy'
'lack'
```

All backslash escape sequences are available for use within a bracket expression. For example, the following expression matches either a caret or a backslash:

```
S LIKE_REGEX '[\\^\\]'
```

The following expression matches all letters or digits:

```
S LIKE_REGEX '[\p{L}\p{Nd}]'
```

where “\p{L}” is the escape matching any letter and “\p{Nd}” is the escape matching any digit.

4.8.3 Matching a range

A minus sign “-” is used to specify a character range. For example:

```
S LIKE_REGEX '[sa-my]'
```

matches the lowercase letters “s”, all the letters between “a” and “m” inclusive, and “y”. Ranges are defined in terms of the UCS code point ordering. When there are multiple ranges, the bracket expression matches the union of the ranges. For example:

```
S LIKE_REGEX '[a-me-z]'
```

matches all lowercase letters.

Using a special character in a range is sometimes permitted, but tricky. Rather than present the rules here, one should use a backslash escape if the start or end point of a range is to be a special character.

4.8.4 Negation

A caret “^” is a special character when it is the first character of a bracket expression, where it indicates that the set of characters is anything not listed by the following bracket expression. For example:

```
S LIKE_REGEX '[^aj-m]'
```

is *True* if *S* contains any character that is not “a”, “j”, “k”, “l”, or “m”.

4.8.5 Character class subtraction

A bracket expression may conclude with a minus sign “-” followed by a nested bracket expression. This is called a *character class subtraction*, and indicates that any character matched by the nested bracket expression is to be removed from the set of characters that might be a match. For example:

```
S LIKE_REGEX '[a-z-[m-p]]'
```

matches anything between “a” and “z”, except for the letters between “m” and “p”, inclusive. This example is equivalent to:

```
S LIKE_REGEX '[a-lq-z]'
```

Character class subtractions can be nested indefinitely (although implementation-defined limitations may apply).

4.9 Alternation

A choice of regular expressions is specified through the use of a vertical bar: “|”. For example:

```
S LIKE_REGEX 'a|b'
```

is *True* if *S* contains either an “a” or a “b”.

Alternation has lower precedence than concatenation. Thus

4.9 Alternation

```
S LIKE_REGEX 'ab|xyz'
```

is *True* if *S* contains either “ab” or “xyz”. This precedence can be overridden through the use of parentheses, such as this example:

```
S LIKE_REGEX 'a(b|xy)z'
```

The preceding example is *True* if *S* contains either “abz” or “axyz”.

4.10 Quantifiers

Quantifiers are metacharacters that specify a match for some number of repetitions of a regular expression. There are two sets of quantifiers, the greedy and the reluctant. The *greedy quantifiers* are:

| | |
|-------|--|
| {n} | Exactly n repetitions. |
| {n,} | n or more repetitions. |
| {n,m} | Between n and m repetitions, inclusive. |
| ? | 0 (zero) or 1 (one) repetition; equivalent to {0,1}. |
| * | 0 (zero) or more repetitions; equivalent to {0,}. |
| + | 1 (one) or more repetitions; equivalent to {1,}. |

The *reluctant quantifiers* are formed by suffixing a question mark to a greedy quantifier. Thus, “*?” is the reluctant form of “*”, and “??” is the reluctant form of “?”. The greedy quantifiers try to match as much as possible, whereas the reluctant quantifiers try to match as little as possible (while still allowing the overall regular expression to match). There is no difference in behavior between the greedy and reluctant quantifiers for LIKE_REGEX. This difference for the other operators is discussed in subsequent Subclauses.

Example:

```
S LIKE_REGEX 'a{3}'
```

is equivalent to

```
S LIKE_REGEX 'aaa'
```

and matches any string containing at least three consecutive instances of “a”. Note that if *S* contains more than three consecutive instances of “a”, it still matches; to test whether *S* contains a substring of three consecutive instances of “a” and no more is a lot harder, since the expression also has to require something other than an “a” at both ends of the substring.

Another example:

```
S LIKE_REGEX 'ab+c'
```

is equivalent to

```
S LIKE_REGEX 'ab{1,}c'
```

and matches any string that contains a substring consisting of an “a”, one or more “b”s, and then a “c”.

4.11 Locating a match

LIKE_REGEX only cares whether a match exists; the other operators care about where a match is located in the string. Consider the regular expression “a+” and the string “a1aa2aaa3”. There are ten possible matches for “a+”, indicated by the underlining on the following lines:

```
'a1aa2aaa3' -- position 1, length 1
'a1aa2aaa3' -- position 3, length 1
'alaa2aaa3' -- position 3, length 2
'alaa2aaa3' -- position 4, length 1

'alaa2aaa3' -- position 6, length 1
'alaa2aaa3' -- position 6, length 2
'alaa2aaa3' -- position 6, length 3
'alaa2aaa3' -- position 7, length 1
'alaa2aaa3' -- position 7, length 2
'alaa2aaa3' -- position 8, length 1
```

Notice that some of the matches are substrings of other matches. The rules of XQuery regular expressions are designed to ignore certain matches, so that the recognized matches are mutually disjoint. Obviously there are many ways to do this, so the rules provide priorities in determining the recognized matches. There are three priorities:

- 1) The top priority is to find a match as early in the string as possible. This is commonly called the *leftmost rule*.
- 2) The second priority is to find the first alternative of an alternation, if possible. There does not appear to be a common name for this rule.
- 3) The last priority is to find the longest possible match for greedy quantifiers, and the shortest match for reluctant quantifiers. In the case of greedy quantifiers, this is commonly called the *longest rule*; there does not appear to be a common name for the rule regarding reluctant quantifiers.

[Historical note: POSIX only has a leftmost longest rule. There were no reluctant quantifiers, and the priority for matching alternations was the longest match rather than the first alternative.]

These rules are illustrated by examples in Table 1, “Match priorities”.

Table 1 — Match priorities

| Subject string | regular expression | match(es) underlined | priority |
|----------------|--------------------|--|---|
| baaaaaa | ba a* | <u>baaaaaa</u> baaaaaa | first alternative; second match always starts after the first match (even though aaaaaa would be longer) |
| ab | a ab | <u>a</u> b | first alternative (rather than matching ab) |
| abcabbabc | ab* | <u>abc</u> abbabc abc <u>ab</u> babc abcabb <u>a</u> b | leftmost longest (greedy quantifier consumes two 'b's) longest |
| abcabbabc | ab*? | <u>a</u> bcabbabc abc <u>a</u> bbabc abcabb <u>a</u> b | shortest (no need to match 'b') shortest shortest |

4.12 Capture and back-reference

A *parenthesized sub-expression* is a portion of a regular expression that is enclosed in parentheses. Parenthesized sub-expressions are numbered in order of their left parenthesis. For example, in the regular expression

```
((a)|(b))
```

there are three sub-expressions:

- 1) ((a)|(b))
- 2) (a)
- 3) (b)

A sub-expression can be referenced later in a regular expression using a back-reference, taking the form of a backslash followed by one or more digits. Thus the three sub-expressions in the example can be referenced as “\1”, “\2”, and “\3”. For example, consider the regular expression:

```
\p{Z}(\p{L}*)\p{Z}*\1\p{Z}
```

The first and only parenthesized sub-expression (“\p{L}”*) matches any sequence of letters that is bounded by some kind of space character (“\p{Z}”) before and after the sequence of letters. The back-reference (“\1”) matches whatever sequence of letters was captured by the first sub-expression. This regular expression might be used to search for occurrences of a repeated word (perhaps caused by a cut-and-paste error). Here is an example of a subject string, with underlining to indicate the match for the entire regular expression:

```
Hello Dolly you're looking looking swell
```

When a back-reference references a parenthesized group with a quantifier, then the back-reference matches the last iteration of the quantified sub-expression. For example, consider the regular expression:

```
'(ab*)*c*\1'
```

and the subject string:

```
'abbbabbabcabbbbb'
```

The matches to “(ab*)” are shown by underlining:

```
'abbbabbabcabbbbb'  
'abbbababbcabbbbb'  
'abbbabbabcabbbbb'
```

These three iterations of “(ab*)” are matched by “(ab*)*” and then the “c” is matched. Next, a match for “\1” is needed. The last match for the first parenthesized sub-expression is “ab”, so the overall match is indicate by underlining:

```
'abbbabbabcabbbbb'
```

In the event that a sub-expression is unmatched, a back-reference to it matches the zero-length string. For example, consider the regular expression:

```
'((a*)|(b*))c??\3'
```

and the subject string:

```
'xyzaaccb'
```

In this example, the alternation “ $((a^*)|(b^*))$ ” matches the “aa”, which is a match for the first alternative. Thus there is no match for the second alternative, “ (b^*) ”. The “ $c?$ ” prefers to match a zero-length string (though it could match the “c”), and the “ \backslash_3 ” in this example always matches a zero-length string. Thus, the complete substring that is matched is underlined:

'xyzaaaccb'

4.13 Precedence

The precedence of operators outside bracket expressions is as follows (from highest to lowest):

— Highest precedence: atoms, defined as:

- Parentheses.
- Individual characters.
- Escape sequences.
- Dot (“.”)
- Anchors (“^”, “\$”)
- Bracket expressions.

— Quantifiers.

— Concatenation.

— Alternation (lowest).

Examples:

1) Quantifiers have higher precedence than concatenation:

ab^* is equivalent to $a(b^*)$

2) Concatenation outranks alternation:

$ab|cd$ is equivalent to $(ab)|(cd)$

4.14 Modes

The preceding discussion has mentioned two of the flags, “s” to specify dot-all mode, and “m” to specify multi-line mode. There are two additional flags, “i” for case-insensitive mode, and “x” to disregard whitespace in regular expressions for readability. The complete set of modes is:

“s” Specifies dot-all mode, in which a period matches any character. If “s” is not specified, then a period matches any single character except a line terminator.

“m” Specifies multi-line mode, in which the anchors match the beginning or end of a line. If “m” is not specified, then the anchors match the beginning or end of the subject string.

“i” Specifies case-insensitive mode.

"x" Specifies that whitespace characters in a regular expression are ignored. This allows application writers to set off portions of a regular expression for greater readability.

IECNORM.COM : Click to view the full PDF of ISO/IEC 19075-1:2021

5 Operators using regular expressions

5.1 Introduction to operators using regular expressions

SQL contains five operators that use the XQuery regular expression syntax:

- 1) `LIKE_REGEX` — predicate that returns *True* if a substring of a string matches a regular expression.
- 2) `OCCURRENCES_REGEX` — numeric function returning the number of matches for a regular expression in a string.
- 3) `POSITION_REGEX` — numeric function returning the position of the start of a match for a regular expression in a string, or the position of the next character after a match.
- 4) `SUBSTRING_REGEX` — character string function returning a substring that matches a regular expression in a string.
- 5) `TRANSLATE_REGEX` — character function that performs a replacement operation on one or all matches to a regular expression in a string.

5.2 LIKE_REGEX

`LIKE_REGEX` is a predicate that returns *True* if a substring of a string matches a regular expression.

The syntax is:

```
<regex like predicate> ::=
  <row value predicand>
  [ NOT ] LIKE_REGEX <XQuery pattern>
  [ FLAG <XQuery option flag> ]
```

where

- `<row value predicand>` is the subject string to be searched for matches to the `<XQuery pattern>`.
- `<XQuery pattern>` is a character string expression whose value is an XQuery regular expression.
- `<XQuery option flag>` is an optional character string, corresponding to the `$flags` argument of the [XQuery and XPath Functions and Operators 3.1](#) function `fn:match`.

The result is *Unknown* if any of the operands is the null value, *True* if there is a substring that matches the `<XQuery pattern>` in the `<row value predicand>`, and *False* if there is no match.

Note that unlike `LIKE`, `LIKE_REGEX` can return *True* without matching the entire string. The usual convention for regular expression matching is to search for a match somewhere within the searched string, without necessarily matching the entire string. The user may use anchors to require a match to the entire string.

Exceptional cases:

- If any of the parameters is the null value, the result is *Unknown*.
- If the pattern or flag is not valid, then an exception condition is raised.

Examples:

'abcde' LIKE_REGEX 'c' evaluates to *True*.

'abcde' LIKE_REGEX 'x' evaluates to *False*.

'abcde' LIKE_REGEX CAST (NULL AS CHAR(10)) evaluates to *Unknown*.

'abcde' LIKE_REGEX '\ ' raises an exception condition. In this example, “\” is not a well-formed regular expression.

'abcde' LIKE_REGEX 'x' FLAG '?' raises an exception condition. In this example, the flag “?” is invalid.

5.3 OCCURRENCES_REGEX

OCCURRENCES_REGEX is a numeric function returning the number of matches for a regular expression in a string. The syntax is:

```
<regex occurrences function> ::=  
  OCCURRENCES_REGEX <left paren>  
    <XQuery pattern> [ FLAG <XQuery option flag> ]  
    IN <regex subject string>  
    [ FROM <start position> ]  
    [ USING <char length units> ] <right paren>
```

where:

- <XQuery pattern> is a character string expression whose value is an XQuery regular expression.
- <XQuery option flag> is an optional character string, corresponding to the \$flags argument of the [XQuery and XPath Functions and Operators 3.1](#) function fn:match.
- <regex subject string> is the character string to be searched for matches to the <XQuery pattern>.
- <start position> is an optional exact numeric value with scale 0 (zero) specifying the position at which to start the search (the default is position 1 (one)).
- <char length units> is CHARACTERS or OCTETS, indicating the unit in which <start position> is measured (the default is to measure in CHARACTERS).

The <regex subject string> is searched for matches to the <XQuery pattern>, starting from position <start position>, which is measured in the units specified by <char length units>, either CHARACTERS or OCTETS. The result is the number of matches.

Exceptional cases:

- If any of the parameters is the null value, then the result is the null value.
- If the pattern or flag is not valid, then an exception condition is raised.
- If a starting position is given in octets, but it is not the first octet of a character, then the result is implementation-dependent. The use of OCTETS is discussed under POSITION_REGEX.
- If any of the numeric parameters is too large or too small, then the result is -1. This includes the following cases:
 - The starting position is less 1 (one).
 - The starting position is greater than the length of the string (measured in CHARACTERS or OCTETS as specified by <char length units>).

Examples:

```
OCCURRENCES_REGEX ( 'a' IN 'what is that?' ) evaluates to 2.
OCCURRENCES_REGEX ( 'a' IN 'what is that?' FROM 5) evaluates to 1 (one).
OCCURRENCES_REGEX ( 'A' FLAG 'i' IN 'what is that' ) evaluates to 2.
OCCURRENCES_REGEX ( 'A' IN 'what is that' ) evaluates to 0 (zero).
```

5.4 POSITION_REGEX

POSITION_REGEX is a numeric function returning the position of the start of a match, or one plus the end of a match, for a regular expression in a string. The syntax is:

```
<regex position expression> ::=
  POSITION_REGEX <left paren> [ <regex position start or after> ]
  <XQuery pattern> [ FLAG <XQuery option flag> ]
  IN <regex subject string>
  [ FROM <start position> ]
  [ USING <char length units> ]
  [ OCCURRENCE <regex occurrence> ]
  [ GROUP <regex capture group> ] <right paren>

<regex position start or after> ::=
  START
  | AFTER
```

where:

- START indicates that the starting position of the match to the regular expression is desired; AFTER indicates that the character position immediately following the match is desired (START is the default). If the match consumes the last character of the subject string, then AFTER returns the length of the string plus 1 (one).
- <XQuery pattern> is a character string expression whose value is an XQuery regular expression.
- <XQuery option flag> is an optional character string, corresponding to the \$flags argument of the [XQuery and XPath Functions and Operators 3.1](#) function `fn:match`.
- <regex subject string> is the character string to be searched for matches to the <XQuery pattern>.
- <start position> is an optional exact numeric value with scale 0 (zero), identifying the character position at which to start the search (the default is 1 (one)).
- <char length units> is CHARACTERS or OCTETS, indicating the unit in which <start position> is measured, and the unit in which the returned position is measured (the default is to measure in CHARACTERS).
- <regex occurrence> is an optional exact numeric value with scale 0 (zero) indicating which occurrence of a match is desired (the default is 1 (one)).
- <regex capture group> is an optional exact numeric value with scale 0 (zero) indicating which capture group of a match is desired (the default is 0 (zero), indicating the entire occurrence).

The <regex subject string> is searched for matches to the <XQuery pattern>. If there are at least *RO* matches, where *RO* is the value of <regex occurrence>, then either the starting position of the *RO*-th match, or the position immediately following the *RO*-th match, is returned (for the START or AFTER options, respectively). Positions are measured in the units specified by <char length units>, either

ISO/IEC 19075-1:2021(E)
5.4 POSITION_REGEX

CHARACTERS or OCTETS. If a <regex capture group> *CAP* is specified, then the position at the start or immediately following the substring that matches the *CAP*-th parenthesized subexpression is used.

With AFTER, note that the position returned is the one after the match. If the match consumes the last character of the string, then the position returned is actually one plus the length of the string (in characters or octets, as requested by <char length units>). The rationale for providing the position that is 1 (one) after the end of the match is that this is the correct place to begin a search for the next match. If the user wishes to process the subject string in a loop, the loop can continue until the AFTER position is greater than the length of the subject string. However, when doing this, the user must beware of a pitfall: if the regular expression matches a zero-length string, then the AFTER position and the START position are the same, and resuming the search at the AFTER position will simply find the same zero-length match again.

OCTETS is provided for efficient processing for those UCS encodings that do not have a fixed character width. It is expected that the user will use the output of POSITION_REGEX (... USING OCTETS ...) to learn the position of some occurrence within a string, measured in octets. That value is then known to be the first octet of a character, and may be used as a starting position in other function invocations. If the user picks an arbitrary octet number, it may be other than the first octet of a character. Naturally, beginning a regular expression match at such an octet can produce unpredictable results. Therefore, the result is said to be implementation-dependent if a starting octet is not the first octet of a character.

Exceptional cases:

- If any of the parameters is the null value, the result is the null value.
- If the pattern or flag is not valid, then an exception condition is raised.
- If a starting position is given in octets, but it is not the first octet of a character, then the result is implementation-dependent.
- If any of the numeric parameters is too large or too small, then the result is 0 (zero). This includes the following cases:
 - The starting position is less 1 (one).
 - The starting position is greater than the length of the string (measured in CHARACTERS or OCTETS as specified by <char length units>).
 - There are not at least *RO* matches.
 - There are no *CAP* parenthesized subexpressions.

Examples:

```
POSITION_REGEX ( 'a' IN 'what is that?' ) evaluates to 3.  
POSITION_REGEX ( START 'a' IN 'what is that?' ) evaluates to 3.  
POSITION_REGEX ( AFTER 'a' IN 'what is that?' ) evaluates to 4.  
POSITION_REGEX ( AFTER 'a' IN 'a' ) evaluates to 2.  
POSITION_REGEX ( 'a' IN 'what is that?' FROM 5 ) evaluates to 11.  
POSITION_REGEX ( 'a' IN 'what is that?' OCCURRENCE 2 ) evaluates to 11.  
POSITION_REGEX ( '(a)(t)' IN 'what is that?' GROUP 2 ) evaluates to 4.  
POSITION_REGEX ( 'A' FLAG 'i' IN 'what is that' ) evaluates to 3.  
POSITION_REGEX ( 'A' IN 'what is that' ) evaluates to 0.
```

5.5 SUBSTRING_REGEX

SUBSTRING_REGEX is a character string function returning a substring that matches a regular expression in a string. The syntax is:

```
<regex substring function> ::=
  SUBSTRING_REGEX <left paren>
    <XQuery pattern> [ FLAG <XQuery option flag> ]
    IN <regex subject string>
    [ FROM <start position> ]
    [ USING <char length units> ]
    [ OCCURRENCE <regex occurrence> ]
    [ GROUP <regex capture group> ] <right paren>
```

where:

- <XQuery pattern> is a character string expression whose value is an XQuery regular expression.
- <XQuery option flag> is an optional character string, corresponding to the \$flags argument of the XQuery and XPath Functions and Operators 3.1 function fn:match.
- <regex subject string> is the character string to be searched for matches to the <XQuery pattern>.
- <start position> is an optional exact numeric value with scale 0 (zero), indicating the character position at which to start the search (the default is position 1 (one)).
- <char length units> is CHARACTERS or OCTETS, indicating the unit in which <start position> is measured (the default is to measure in CHARACTERS).
- <regex occurrence> is an optional exact numeric value with scale 0 (zero) indicating which occurrence of a match is desired (the default is 1 (one)).
- <regex capture group> is an optional exact numeric value with scale 0 (zero) indicating which capture group of a match is desired (the default is 0 (zero), indicating the entire occurrence).

The <regex subject string> is searched for matches to the <XQuery pattern>. If there are at least *RO* matches, where *RO* is the value of <regex occurrence>, then the result is the substring that is the *RO*-th match. If <regex capture group> *CAP* is specified, then the result is the substring that matches the *CAP*-th parenthesized substring within the substring that is the *RO*-th match. If there are not at least *RO* matches, or at least *CAP* parenthesized subexpressions, the result is the null value.

The exceptional cases are:

- If any of the parameters is the null value, the result is the null value.
- If the pattern or flag is not valid, then an exception condition is raised.
- If a starting position is given in octets, but it is not the first octet of a character, then the result is implementation-dependent.
- If any of the numeric parameters is too large or too small, then the result is the null value. This includes the following cases:
 - The starting position is less than 1 (one).
 - The starting position is greater than the length of the string (measured in CHARACTERS or OCTETS as specified by <char length units>).
 - There are not at least *RO* matches.
 - There are not *CAP* parenthesized subexpressions.