

2.3. pyspark.sql.functions

A collections of builtin functions

2.3.1. Functions

<code>abs</code> (col)	Computes the absolute value.
<code>acos</code> (col)	Computes the cosine inverse of the given value; the returned
<code>add_months</code> (start, months)	Returns the date that is <i>months</i> months after <i>start</i>
<code>approxCountDistinct</code> (col[, rsd])	Returns a new <code>column</code> for approximate distinct count of <code>col</code>
<code>array</code> (*cols)	Creates a new array column.
<code>array_contains</code> (col, value)	Collection function: returns True if the array contains the giv
<code>asc</code> (col)	Returns a sort expression based on the ascending order of th
<code>ascii</code> (col)	Computes the numeric value of the first character of the stri
<code>asin</code> (col)	Computes the sine inverse of the given value; the returned a
<code>atan</code> (col)	Computes the tangent inverse of the given value.
<code>atan2</code> (col1, col2)	Returns the angle theta from the conversion of rectangular c
<code>avg</code> (col)	Aggregate function: returns the average of the values in a gr
<code>base64</code> (col)	Computes the BASE64 encoding of a binary column and retu
<code>bin</code> (col)	Returns the string representation of the binary value of the g

<code>bitwiseNOT</code> (col)	Computes bitwise not.
<code>broadcast</code> (df)	Marks a DataFrame as small enough for use in broadcast join
<code>bround</code> (col[, scale])	Round the given value to <i>scale</i> decimal places using HALF_UP
<code>cbrt</code> (col)	Computes the cube-root of the given value.
<code>ceil</code> (col)	Computes the ceiling of the given value.
<code>coalesce</code> (*cols)	Returns the first column that is not null.
<code>col</code> (col)	Returns a <code>column</code> based on the given column name.
<code>collect_list</code> (col)	Aggregate function: returns a list of objects with duplicates.
<code>collect_set</code> (col)	Aggregate function: returns a set of objects with duplicate elements
<code>column</code> (col)	Returns a <code>column</code> based on the given column name.
<code>concat</code> (*cols)	Concatenates multiple input string columns together into a single string
<code>concat_ws</code> (sep, *cols)	Concatenates multiple input string columns together into a single string with a separator
<code>conv</code> (col, fromBase, toBase)	Convert a number in a string column from one base to another
<code>corr</code> (col1, col2)	Returns a new <code>column</code> for the Pearson Correlation Coefficient
<code>cos</code> (col)	Computes the cosine of the given value.
<code>cosh</code> (col)	Computes the hyperbolic cosine of the given value.
<code>count</code> (col)	Aggregate function: returns the number of items in a group.
<code>countDistinct</code> (col, *cols)	Returns a new <code>column</code> for distinct count of <code>col</code> or <code>cols</code> .
<code>covar_pop</code> (col1, col2)	Returns a new <code>column</code> for the population covariance of <code>col1</code> and <code>col2</code>
<code>covar_samp</code> (col1, col2)	Returns a new <code>column</code> for the sample covariance of <code>col1</code> and <code>col2</code>
<code>crc32</code> (col)	Calculates the cyclic redundancy check value (CRC32) of a binary value
<code>create_map</code> (*cols)	Creates a new map column.

<code>cume_dist</code> ()	Window function: returns the cumulative distribution of values.
<code>current_date</code> ()	Returns the current date as a date column.
<code>current_timestamp</code> ()	Returns the current timestamp as a timestamp column.
<code>date_add</code> (start, days)	Returns the date that is <i>days</i> days after <i>start</i> .
<code>date_format</code> (date, format)	Converts a date/timestamp/string to a value of string in the format.
<code>date_sub</code> (start, days)	Returns the date that is <i>days</i> days before <i>start</i> .
<code>datediff</code> (end, start)	Returns the number of days from <i>start</i> to <i>end</i> .
<code>dayofmonth</code> (col)	Extract the day of the month of a given date as integer.
<code>dayofyear</code> (col)	Extract the day of the year of a given date as integer.
<code>decode</code> (col, charset)	Computes the first argument into a string from a binary using the charset.
<code>dense_rank</code> ()	Window function: returns the rank of rows within a window.
<code>desc</code> (col)	Returns a sort expression based on the descending order of the column.
<code>encode</code> (col, charset)	Computes the first argument into a binary from a string using the charset.
<code>exp</code> (col)	Computes the exponential of the given value.
<code>explode</code> (col)	Returns a new row for each element in the given array or map.
<code>expm1</code> (col)	Computes the exponential of the given value minus one.
<code>expr</code> (str)	Parses the expression string into the column that it represents.
<code>factorial</code> (col)	Computes the factorial of the given value.
<code>first</code> (col[, ignorenulls])	Aggregate function: returns the first value in a group.
<code>floor</code> (col)	Computes the floor of the given value.
<code>format_number</code> (col, d)	Formats the number X to a format like '#,-#,-#.-', rounded to d decimal places.
<code>format_string</code> (format, *cols)	Formats the arguments in printf-style and returns the result.

<code>from_unixtime</code> (timestamp[, format])	Converts the number of seconds from unix epoch (1970-01-01 00:00:00 UTC) to a timestamp.
<code>from_utc_timestamp</code> (timestamp, tz)	Assumes given timestamp is UTC and converts to given time zone.
<code>get_json_object</code> (col, path)	Extracts json object from a json string based on json path specification.
<code>greatest</code> (*cols)	Returns the greatest value of the list of column names, skipping null values in the list.
<code>grouping</code> (col)	Aggregate function: indicates whether a specified column in the group by clause is part of the grouping.
<code>grouping_id</code> (*cols)	Aggregate function: returns the level of grouping, equals to the sum of the grouping function.
<code>hash</code> (*cols)	Calculates the hash code of given columns, and returns the result as a long.
<code>hex</code> (col)	Computes hex value of the given column, which could be String or binary.
<code>hour</code> (col)	Extract the hours of a given date as integer.
<code>hypot</code> (col1, col2)	Computes $\sqrt{a^2 + b^2}$ without intermediate overflow or underflow.
<code>ignore_unicode_prefix</code> (f)	Ignore the 'u' prefix of string in doc tests, to make it works with older versions of Spark.
<code>initcap</code> (col)	Translate the first letter of each word to upper case in the separator specified.
<code>input_file_name</code> ()	Creates a string column for the file name of the current Spark task.
<code>instr</code> (str, substr)	Locate the position of the first occurrence of substr column in str.
<code>isnan</code> (col)	An expression that returns true iff the column is NaN.
<code>isnull</code> (col)	An expression that returns true iff the column is null.
<code>json_tuple</code> (col, *fields)	Creates a new row for a json column according to the given fields.
<code>kurtosis</code> (col)	Aggregate function: returns the kurtosis of the values in a group.
<code>lag</code> (col[, count, default])	Window function: returns the value that is <i>offset</i> rows before the current row.
<code>last</code> (col[, ignorenulls])	Aggregate function: returns the last value in a group.
<code>last_day</code> (date)	Returns the last day of the month which the given date belongs to.
<code>lead</code> (col[, count, default])	Window function: returns the value that is <i>offset</i> rows after the current row.

<code>least</code> (*cols)	Returns the least value of the list of column names, skipping nulls.
<code>length</code> (col)	Calculates the length of a string or binary expression.
<code>levenshtein</code> (left, right)	Computes the Levenshtein distance of the two given strings.
<code>lit</code> (col)	Creates a <code>column</code> of literal value.
<code>locate</code> (substr, str[, pos])	Locate the position of the first occurrence of substr in a string.
<code>log</code> (arg1[, arg2])	Returns the first argument-based logarithm of the second argument.
<code>log10</code> (col)	Computes the logarithm of the given value in Base 10.
<code>log1p</code> (col)	Computes the natural logarithm of the given value plus one.
<code>log2</code> (col)	Returns the base-2 logarithm of the argument.
<code>lower</code> (col)	Converts a string column to lower case.
<code>lpad</code> (col, len, pad)	Left-pad the string column to width <i>len</i> with <i>pad</i> .
<code>ltrim</code> (col)	Trim the spaces from left end for the specified string value.
<code>max</code> (col)	Aggregate function: returns the maximum value of the expression.
<code>md5</code> (col)	Calculates the MD5 digest and returns the value as a 32 character string.
<code>mean</code> (col)	Aggregate function: returns the average of the values in a group.
<code>min</code> (col)	Aggregate function: returns the minimum value of the expression.
<code>minute</code> (col)	Extract the minutes of a given date as integer.
<code>monotonically_increasing_id</code> ()	A column that generates monotonically increasing 64-bit integers.
<code>month</code> (col)	Extract the month of a given date as integer.
<code>months_between</code> (date1, date2)	Returns the number of months between date1 and date2.
<code>nanvl</code> (col1, col2)	Returns col1 if it is not NaN, or col2 if col1 is NaN.
<code>next_day</code> (date, dayOfWeek)	Returns the first date which is later than the value of the date argument.

<code>ntile</code> (n)	Window function: returns the ntile group id (from 1 to n incl
<code>percent_rank</code> ()	Window function: returns the relative rank (i.e.
<code>posexplode</code> (col)	Returns a new row for each element with position in the give
<code>pow</code> (col1, col2)	Returns the value of the first argument raised to the power o
<code>quarter</code> (col)	Extract the quarter of a given date as integer.
<code>rand</code> ([seed])	Generates a random column with i.i.d.
<code>randn</code> ([seed])	Generates a column with i.i.d.
<code>rank</code> ()	Window function: returns the rank of rows within a window
<code>regexp_extract</code> (str, pattern, idx)	Extract a specific(idx) group identified by a java regex, from t
<code>regexp_replace</code> (str, pattern, replacement)	Replace all substrings of the specified string value that match
<code>repeat</code> (col, n)	Repeats a string column n times, and returns it as a new strin
<code>reverse</code> (col)	Reverses the string column and returns it as a new string col
<code>rint</code> (col)	Returns the double value that is closest in value to the argum
<code>round</code> (col[, scale])	Round the given value to <i>scale</i> decimal places using HALF_U
<code>row_number</code> ()	Window function: returns a sequential number starting at 1,
<code>rpad</code> (col, len, pad)	Right-pad the string column to width <i>len</i> with <i>pad</i> .
<code>rtrim</code> (col)	Trim the spaces from right end for the specified string value.
<code>second</code> (col)	Extract the seconds of a given date as integer.
<code>sha1</code> (col)	Returns the hex string result of SHA-1.
<code>sha2</code> (col, numBits)	Returns the hex string result of SHA-2 family of hash functio
<code>shiftLeft</code> (col, numBits)	Shift the given value numBits left.
<code>shiftRight</code> (col, numBits)	Shift the given value numBits right.

<code>shiftRightUnsigned</code> (col, numBits)	Unsigned shift the given value numBits right.
<code>signum</code> (col)	Computes the signum of the given value.
<code>sin</code> (col)	Computes the sine of the given value.
<code>since</code> (version)	A decorator that annotates a function to append the version
<code>sinh</code> (col)	Computes the hyperbolic sine of the given value.
<code>size</code> (col)	Collection function: returns the length of the array or map st
<code>skewness</code> (col)	Aggregate function: returns the skewness of the values in a g
<code>sort_array</code> (col[, asc])	Collection function: sorts the input array for the given colum
<code>soundex</code> (col)	Returns the SoundEx encoding for a string
<code>spark_partition_id</code> ()	A column for partition ID of the Spark task.
<code>split</code> (str, pattern)	Splits str around pattern (pattern is a regular expression).
<code>sqrt</code> (col)	Computes the square root of the specified float value.
<code>stddev</code> (col)	Aggregate function: returns the unbiased sample standard d
<code>stddev_pop</code> (col)	Aggregate function: returns population standard deviation o
<code>stddev_samp</code> (col)	Aggregate function: returns the unbiased sample standard d
<code>struct</code> (*cols)	Creates a new struct column.
<code>substring</code> (str, pos, len)	Substring starts at <i>pos</i> and is of length <i>len</i> when str is String t
<code>substring_index</code> (str, delim, count)	Returns the substring from string str before count occurrenc
<code>sum</code> (col)	Aggregate function: returns the sum of all values in the expres
<code>sumDistinct</code> (col)	Aggregate function: returns the sum of distinct values in the
<code>tan</code> (col)	Computes the tangent of the given value.
<code>tanh</code> (col)	Computes the hyperbolic tangent of the given value.

<code>toDegrees</code> (col)	Converts an angle measured in radians to an approximately ε
<code>toRadians</code> (col)	Converts an angle measured in degrees to an approximately
<code>to_date</code> (col)	Converts the column of StringType or TimestampType into D
<code>to_utc_timestamp</code> (timestamp, tz)	Assumes given timestamp is in given timezone and converts t
<code>translate</code> (srcCol, matching, replace)	A function translate any character in the <i>srcCol</i> by a character
<code>trim</code> (col)	Trim the spaces from both ends for the specified string column
<code>trunc</code> (date, format)	Returns date truncated to the unit specified by the format.
<code>udf</code> (f[, returnType])	Creates a <code>column</code> expression representing a user defined fun
<code>unbase64</code> (col)	Decodes a BASE64 encoded string column and returns it as a
<code>unhex</code> (col)	Inverse of hex.
<code>unix_timestamp</code> ([timestamp, format])	Convert time string with given pattern ('yyyy-MM-dd HH:mr
<code>upper</code> (col)	Converts a string column to upper case.
<code>v</code> (name[, doc])	Create a binary mathfunction by name
<code>var_pop</code> (col)	Aggregate function: returns the population variance of the v
<code>var_samp</code> (col)	Aggregate function: returns the unbiased variance of the val
<code>variance</code> (col)	Aggregate function: returns the population variance of the v
<code>weekofyear</code> (col)	Extract the week number of a given date as integer.
<code>when</code> (condition, value)	Evaluates a list of conditions and returns one of multiple pos:
<code>window</code> (timeColumn, windowDuration[, ...])	Bucketize rows into one or more time windows given a times
<code>year</code> (col)	Extract the year of a given date as integer.

2.3.2. Classes

<code>AutoBatchedSerializer</code> (serializer[, bestSize])	Choose the size of batch automatically based on the siz
<code>Column</code> (jc)	A column in a DataFrame.
<code>DataFrame</code> (jdf, sql_ctx)	A distributed collection of data grouped into named col
<code>PickleSerializer</code> ()	Serializes objects using Python's pickle serializer:
<code>SparkContext</code> ([master, appName, sparkHome, ...])	Main entry point for Spark functionality.
<code>StringType</code>	String data type.
<code>UserDefinedFunction</code> (func, returnType[, name])	User defined function in Python