

We Do Hadoop



Cheat Sheet **Hive Functions**

Contents

- 1 Creating Functions
- 2 Mathematical, String, Date, Collection, Text, Conditional
- Built-In Aggregates, Built-In Table-Generating Functions

HiveQL provides big data querying capabilities in Hadoop in a dialect very similar to SQL and should be familiar to anyone used to working with SQL databases. This cheat sheet covers more advanced capabilities of Hive – specifically creating and using User Defined Functions (UDFs) in Hive, and some of the built in functions

Additional Resources



Learn to become fluent in Apache Hive with the Hive Language Manual: https://cwiki.apache.org/confluence/display/Hive/LanguageManual



Try out Qubole Data Service (QDS) for free: http://www.qubole.com/try



Get in the Hortonworks Sandbox and try out Hadoop with interactive tutorials: http://hortonworks.com/sandbox





Hive Function Meta Commands

 SHOW FUNCTIONS
 Lists Hive functions and operators

 DESCRIBE FUNCTION [function name]
 Displays short description of the function

 DESCRIBE FUNCTION EXTENDED [function name]
 Access extended description of the function

Types of Hive Functions

UDF A function that takes one or more columns from a row as argument and returns a single value or object

e.g. concat(col1, col2)

UDAF Aggregates column values in multiple rows and returns a single value e.g. sum(c1)

UDTF Takes zero or more inputs and produces multiple columns or rows of output e.g. explode()

Macros A function that users other Hive functions

How To Develop UDFs

How To Develop UDFs, Generic UDFs, UDAFs and UDTFs

```
public class YourUDFName extends UDF{
public class YourGenericUDFName extends GenericUDF {..}}
public class YourGenericUDAFName extends AbstractGenericUDAFResolver {..}}
public class YourGenericUDTFName extends GenericUDTF {..}
```

How To Deploy/Drop UDFs

```
At start of each session:
```

```
ADD JAR /full_path_to_jar/YourUDFName.jar;
CREATE TEMPORARY FUNCTION YourUDFName AS 'org.apache.hadoop.hive.contrib.udf.example.YourUDFName';
At the end of each session:
```

DROP TEMPORARY FUNCTION IF EXISTS YourUDFName;





Mathematical Functions

Return Type Name (Signature)		Description	
bigint	round(double a)	Returns the rounded BIGINT value of the double	
double	round(double a, int d)	Returns the double rounded to d decimal places	
bigint	floor(double a)	Returns the maximum BIGINT value that is equal or less than the double	
bigint	<pre>ceil(double a), ceiling(double a)</pre>	Returns the minimum BIGINT value that is equal or greater than the double	
double	rand(), rand(int seed)	Returns a random number (that changes from row to row) that is distributed uniformly from 0 to 1. Specifiying the seed will make sure the generated random number sequence is deterministic.	
double	exp(double a)	Returns e where e is the base of the natural logarithm	
double	ln(double a)	Returns the natural logarithm of the argument	
double	log10(double a)	Returns the base-10 logarithm of the argument	
double	log2(double a)	Returns the base-2 logarithm of the argument	
double	log(double base, double a)	Return the base "base" logarithm of the argument	
double	<pre>pow(double a, double p), power(double a, double p)</pre>	Return a	
double	sqrt(double a)	Returns the square root of a	
string	bin(BIGINT a)	Returns the number in binary format	
string	hex(BIGINT a) hex(string a)	If the argument is an int, hex returns the number as a string in hex format. Otherwise if the number is a string, it converts each character into its hex representation and returns the resulting string.	
string	unhex(string a)	Inverse of hex. Interprets each pair of characters as a hexidecimal number and converts to the character represented by the number.	
string	<pre>conv(BIGINT num, int from_base, int to_base), conv(STRING num, int from_base, int to_base)</pre>	Converts a number from a given base to another	
double	abs(double a)	Returns the absolute value	
int double	<pre>pmod(int a, int b) pmod(double a, double b)</pre>	Returns the positive value of a mod b	
double	sin(double a)	Returns the sine of a (a is in radians)	
double	asin(double a)	Returns the arc sin of x if -1<=a<=1 or null otherwise	
double	cos(double a)	Returns the cosine of a (a is in radians)	
double	acos(double a)	Returns the arc cosine of x if -1<=a<=1 or null otherwise	
double	tan(double a)	Returns the tangent of a (a is in radians)	
double	atan(double a)	Returns the arctangent of a	
double	degrees(double a)	Converts value of a from radians to degrees	
double	radians(double a)	Converts value of a from degrees to radians	
int double	<pre>positive(int a), positive(double a)</pre>	Returns a	
int double	<pre>negative(int a), negative(double a)</pre>	Returns -a	

Twitter: twitter.com/hortonworks LinkedIn: linkedin.com/company/hortonworks





String Functions

Return Type	Name (Signature)	Description	
int	ascii(string str)	Returns the numeric value of the first character of str	
string	concat(string binary A, string binary B)	Returns the string or bytes resulting from concatenating the strings or bytes passed in as parameters in order. e.g. concat('foo', 'bar') results in 'foobar'. Note that this function can take any number of input strings.	
string	<pre>concat_ws(string SEP, string A, string B)</pre>	Like concat() above, but with custom separator SEP.	
string	<pre>concat_ws(string SEP, array<string>)</string></pre>	Like concat_ws() above, but taking an array of strings. (as of Hive 0.9.0)	
int	<pre>find_in_set(string str, string strList)</pre>	Returns the first occurance of str in strList where strList is a comma-delimited string. Returns null if either argument is null. Returns 0 if the first argument contains any commas. e.g. find_in_set('ab', 'abc,b,ab,c,def') returns 3	
string	<pre>format_number(number x, int d)</pre>	Formats the number X to a format like '#,###,###, rounded to D decimal places, and returns the result as a string. If D is 0, the result has no decimal point or fractional part. (as of Hive 0.10.0)	
string	<pre>get_json_object(string json_string, string path)</pre>	Extract json object from a json string based on json path specified, and return json string of the extracted json object. It will return null if the input json string is invalid.	
boolean	<pre>in_file(string str, string filename)</pre>	Returns true if the string str appears as an entire line in filename.	
int	instr(string str, string substr)	Returns the position of the first occurence of substr in str	
int	length(string A)	Returns the length of the string	
int	locate(string substr, string str[, int pos])	Returns the position of the first occurrence of substr in str after position pos	
string	lower(string A) lcase(string A)	Returns the string resulting from converting all characters of B to lower case e.g. lower('fOoBaR') results in 'foobar'	
string	lpad(string str, int len, string pad)	Returns str, left-padded with pad to a length of len	
string	ltrim(string A)	Returns the string resulting from trimming spaces from the beginning(left hand side) of A e.g. Itrim('foobar') results in 'foobar'	
string	<pre>parse_url(string urlString, string partToExtract [, string keyToExtract])</pre>	Returns the specified part from the URL. Valid values for partToExtract include HOST, PATH, QUERY, REF, PROTOCOL, AUTHORITY, FILE, and USERINFO.	
string	printf(String format, Obj args)	Returns the input formatted according do printf-style format strings (as of Hive 0.9.0)	
string	regexp_extract(string subject, string pattern, int index)	Returns the string extracted using the pattern. e.g. regexp_extract('foothebar', 'foo(.*?)(bar)', 2) returns 'bar.'	
int double	<pre>pmod(int a, int b) pmod(double a, double b)</pre>	s' is necessary to match whitespace, etc. The 'index' parameter is the Java regex Matcher group() method index. See docs/api/java/util/regex/Matcher.html for more information on the 'index' or Java regex group() method.	
string	regexp_replace(string INITIAL_STRING, string PATTERN, string REPLACEMENT)	Returns the string resulting from replacing all substrings in INITIAL_STRING that match the java regular expression syntax defined in PATTERN with instances of REPLACEMENT, e.g. regexp_replace("foobar", "oo ar", "") returns 'fb.'	
double	asin(double a)	s' is necessary to match whitespace, etc.	
string	repeat(string str, int n)	Repeat str n times	
string	reverse(string A)	Returns the reversed string	
string	rpad(string str, int len, string pad)	Returns str, right-padded with pad to a length of len	
string	rtrim(string A)	Returns the string resulting from trimming spaces from the end(right hand side) of A e.g. rtrim(' foobar ') results in ' foobar'	
array <array<string>></array<string>	sentences(string str, string lang, string locale)	Tokenizes a string of natural language text into words and sentences, where each sentence is broken at the appropriate sentence boundary and returned as an array of words.	
string	space(int n)	Return a string of n spaces	
array	split(string str, string pat)	Split str around pat (pat is a regular expression)	
map <string,string></string,string>	<pre>str_to_map(text[, delimiter1, delimiter2])</pre>	Splits text into key-value pairs using two delimiters. Delimiter1 separates text into K-V pairs, and Delimiter2 splits each K-V pair. Default delimiters are ',' for delimiter1 and '=' for delimiter2.	
string	<pre>substr(string binary A, int start) substring(string binary A, int start)</pre>	Returns the substring or slice of the byte array of A starting from start position till the end of string A e.g. substr('foobar', 4) results in 'bar'	
string	substr(string binary A, int start, int len) substring(string binary A, int start, int len)	Returns the substring or slice of the byte array of A starting from start position with length len e.g. substr('foobar', 4, 1) results in 'b'	
string	translate(string input, string from, string to)	Translates the input string by replacing the characters present in the from string with the corresponding characters in the to string. This is similar to the translatefunction in PostgreSQL. If any of the parameters to this UDF are NULL, the result is NULL as well (available as of Hive 0.10.0)	
string	trim(string A)	Returns the string resulting from trimming spaces from both ends of A e.g. trim(' foobar ') results in 'foobar'	
string	upper(string A) ucase(string A)	Returns the string resulting from converting all characters of A to upper case e.g. upper('fOoBaR') results in 'FOOBAR'	

Twitter: twitter.com/hortonworks LinkedIn: linkedin.com/company/hortonworks





Date Functions

Return Type	Name (Signature)	Description
BIGINT	round(double a)	Returns the rounded BIGINT value of the double
DOUBLE	round(double a, int d)	Returns the double rounded to d decimal places
BIGINT	floor(double a)	Returns the maximum BIGINT value that is equal or less than the double
BIGINT	ceil(double a), ceiling(double a)	Returns the minimum BIGINT value that is equal or greater than the double

Collection Functions

Return Type	Name (Signature)	Description	
int	size(Map <k.v>)</k.v>	Returns the number of elements in the map type	
int	size(Array <t>)</t>	Returns the number of elements in the array type	
array <k></k>	<pre>map_keys(Map<k.v>)</k.v></pre>	Returns an unordered array containing the keys of the input map	
array <v></v>	<pre>map_values(Map<k.v>)</k.v></pre>	Returns an unordered array containing the values of the input map	
boolean array_contains(Array <t>, value)</t>		Returns TRUE if the array contains value	
array <t> sort_array(Array<t>)</t></t>		Sorts the input array in ascending order according to the natural ordering of the array elements and returns it (as of version 0.9.0)	

Text Analytics Functions

Return Type	Name (Signature)	Description
array <struct<string,double>></struct<string,double>	<pre>context_ngrams(array<array<string>>,</array<string></pre>	Returns the top-k contextual N-grams from a set of tokenized sentences, given a string of "context". See StatisticsAndDataMining for more information.
		N-grams are subsequences of length N drawn from a longer sequence. The purpose of the ngrams() UDAF is to find the k most frequent n-grams from one or more sequences. It can be used in conjunction with the sentences() UDF to analyze unstructured natural language text, or the collect() function to analyze more general string data.
array <struct<string,double>></struct<string,double>	ngrams(array <array<string>>, int N, int K, int pf)</array<string>	Returns the top-k N-grams from a set of tokenized sentences, such as those returned by the sentences() UDAF.
		Contextual n-grams are similar to n-grams, but allow you to specify a 'context' string around which n-grams are to be estimated. For example, you can specify that you're only interested in finding the most common two-word phrases in text that follow the context "I love". You could achieve the same result by manually stripping sentences of non-contextual content and then passing them to ngrams(), but context_ngrams() makes it much easier.

Conditional Functions

Return Type	Name (Signature)	Description
Т	<pre>if(boolean testCondition, T valueTrue, T valueFalseOrNull)</pre>	Return valueTrue when testCondition is true, returns valueFalseOrNull otherwise
Т	COALESCE(T v1, T v2,)	Return the first v that is not NULL, or NULL if all v's are NULL
Т	CASE a WHEN b THEN c [WHEN d THEN e]* [ELSE f] END	When a = b, returns c; when a = d, return e; else return f
Т	CASE WHEN a THEN b [WHEN c THEN d]* [ELSE e] END	When a = true, returns b; when c = true, return d; else return e

Twitter: twitter.com/hortonworks LinkedIn: linkedin.com/company/hortonworks

Call: 855-HADOOP-HELP Web: qubole.com/try





Built-In Aggregate Functions (UDAF)

Return Type	Name (Signature)	Description	
bigint	<pre>count(*), count(expr), count(DISTINCT expr[, expr])</pre>	count(*) - Returns the total number of retrieved rows, including rows containing NULL values; count(expr)	
double	sum(col), sum(DISTINCT col)	Returns the sum of the elements in the group or the sum of the distinct values of the column in the group	
double	avg(col), avg(DISTINCT col)	Returns the average of the elements in the group or the average of the distinct values of the column in the group	
double	min(col)	Returns the minimum of the column in the group	
double	max(col)	Returns the maximum value of the column in the group	
double	variance(col), var_pop(col)	Returns the variance of a numeric column in the group	
double	var_samp(col)	Returns the unbiased sample variance of a numeric column in the group	
double	stddev_pop(col)	Returns the standard deviation of a numeric column in the group	
double	stddev_samp(col)	Returns the unbiased sample standard deviation of a numeric column in the group	
double	covar_pop(col1, col2)	Returns the population covariance of a pair of numeric columns in the group	
double	covar_samp(col1, col2)	Returns the sample covariance of a pair of a numeric columns in the group	
double	corr(col1, col2)	Returns the Pearson coefficient of correlation of a pair of a numeric columns in the group	
double	percentile(BIGINT col, p)	Returns the exact p percentile of a column in the group (does not work with floating point types). p must be between 0 and 1.	
array <double></double>			
double	percentile_approx(DOUBLE col, p [, B]) Returns an approximate p percentile of a numeric column (in point types) in the group. The B parameter controls approxima at the cost of memory.		
array <double></double>	<pre>percentile_approx(DOUBLE col, array(p₁ [, p₂]) [, B])</pre>	Same as above, but accepts and returns an array of percentile values instead of a single one.	
array <struct {'x','y'}=""></struct>			
array	collect_set(col)	Returns a set of objects with duplicate elements eliminated	

Built-In Table-Generating Functions (UDTF)

Return Type	Name (Signature)	Description
<pre>inline(ARRAY<struct[,struct]>)</struct[,struct]></pre>	Explode	explode() takes in an array as an input and outputs the elements of the array as separate rows. UDTF's can be used in the SELECT expression list and as a part of LATERAL VIEW.Eg:
		An example use of explode() in the SELECT expression list is as follows: Consider a table named myTable that has a single column (myCol) and two rows: Array <int> myCol [1,2] [4] Then running the query: SELECT explode(myCol) AS myNewCol FROM myTable will produce: (int) myNewCol 1 2</int>

Twitter: twitter.com/hortonworks Linkedin: linkedin.com/company/hortonworks

Call: 855-HADOOP-HELP **Web:** qubole.com/try