

Hive Function Cheat sheet

(<http://quole2.wpengine.com/wp-content/uploads/2014/01/hive-function-cheat-sheet.pdf>)

DOWNLOAD

Hive Function Meta commands

- `SHOW FUNCTION` – 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

Hive Function Cheat sheet

- Date Functions
- Mathematical Functions
- String Functions
- Collection Functions

- UDF— is a function that takes one or more columns from a row as argument and returns a single value or object. g: concat(col1, col2)
- UDTF— takes zero or more inputs and produces multiple columns or rows of output. g: explode()
- Macros— a function that uses other Hive functions.

How To Develop UDFs

package org.apache.hadoop.hive.contrib.udf.example;

```
import java.util.Date;
import java.text.SimpleDateFormat;
import org.apache.hadoop.hive ql.exec.UDF;

@Description(name = "YourUDFName",
    value = "_FUNC_(InputDataType) - using the input data
    type X argument, "+
        "returns YYY.",
    extended = "Example:\n"
        + "    > SELECT _FUNC_(InputDataType) FROM tabl
    ename;")

public class YourUDFName extends UDF{
    ..
    public YourUDFName( InputDataType InputValue ) {
        ..
    }

    public String evaluate( InputDataType InputValue ) {
        ..
    }
}
```

How To Develop UDFs, GenericUDFs, UDAFs, and UDTFs

public class YourUDFName extends UDF{

- UDAF
- UDTF
- Conditional Functions
- Functions for Text Analytics

Go to

Pig Function Cheat sheet

(<http://www.quole.com/resources/cheatsheet/pig-function-cheat-sheet/>)

- 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;
```

Date Functions

The following built-in date functions are supported in hive:

Return Type	Name(signature)	Example
string	from_unixtime(bigint unixtime[, string format])	Converts the number of seconds from unix epoch (1970-01-01 00:00:00 UTC) to a string representing the timestamp of that moment in the current system time zone in the format of “1970-01-01 00:00:00”
bigint	unix_timestamp()	Gets current time stamp using the default time zone.
bigint	unix_timestamp(string date)	Converts time string in format-MM-dd HH:mm:ss to Unix time stamp, return 0 if fail: unix_timestamp('2009-03-20 11:30:01') = 1237573801

igint	unix_timestamp(string date, string pattern)	Convert time string with given pattern to Unix time stamp, return 0 if fail: unix_timestamp('2009-03-20', '-MM-dd') = 1237532400
string	to_date(string timestamp)	Returns the date part of a timestamp string: to_date("1970-01-01 00:00:00") = "1970-01-01"
int	ear(string date)	Returns the ear part of a date or a timestamp string: ear("1970-01-01 00:00:00") = 1970, ear("1970-01-01") = 1970
int	month(string date)	Returns the month part of a date or a timestamp string: month("1970-11-01 00:00:00") = 11, month("1970-11-01") = 11
int	da (string date) da ofmonth(date)	Return the da part of a date or a timestamp string: da ("1970-11-01 00:00:00") = 1, da ("1970-11-01") = 1
int	hour(string date)	Returns the 'hour of the timestamp: hour('2009-07-30 12:58:59') = 12, hour('12:58:59') = 12
int	minute(string date)	Returns the minute of the timestamp
int	second(string date)	Returns the second of the timestamp
int	weekof ear(string date)	Return the week number of a timestamp string: weekof ear("1970-11-01 00:00:00") = 44, weekof ear("1970-11-01") = 44
int	datediff(string enddate, string startdate)	Return the number of days from startdate to enddate: datediff('2009-03-01', '2009-02-27') = 2

string	date_add(string startdate, int days)	Add a number of days to startdate: date_add('2008-12-31', 1) = '2009-01-01'
string	date_sub(string startdate, int days)	Subtract a number of days to startdate: date_sub('2008-12-31', 1) = '2008-12-30'
timestamp	from_utc_timestamp(timestamp, string timezone)	Assumes given timestamp is UTC and converts to given timezone (as of Hive 0.8.0)
timestamp	to_utc_timestamp(timestamp, string timezone)	Assumes given timestamp is in given timezone and converts to UTC (as of Hive 0.8.0)

Mathematical Functions

The following built-in mathematical functions are supported in hive; most return NULL when the argument(s) are NULL:

Return Type	Name(signature)	xample
INTEGER	round(double a)	Returns the rounded INTEGER value of the double
DOUBLE	round(double a, int d)	Returns the double rounded to d decimal places
INTEGER	floor(double a)	Returns the maximum INTEGER value that is equal or less than the double
INTEGER	ceil(double a), ceiling(double a)	Returns the minimum INTEGER 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. Specifying the

		seed will make sure the generated random number sequence is deterministic.
dou le	exp(dou le a)	Returns e^a where e is the base of the natural logarithm
dou le	ln(dou le a)	Returns the natural logarithm of the argument
dou le	log10(dou le a)	Returns the base-10 logarithm of the argument
dou le	log2(dou le a)	Returns the base-2 logarithm of the argument
dou le	log(dou le base, dou le a)	Return the base "base" logarithm of the argument
dou le	pow(dou le a, dou le p), power(dou le a, dou le p)	Return a^p
dou le	sqrt(dou le a)	Returns the square root of a
string	in(INT a)	Returns the number in binary format
string	hex(INT 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 hexadecimal number and converts to the character represented by the number.
string	conv(INT num, int from_base, int to_base), conv(STRING num, int from_base, int to_base)	Converts a number from a given base to another

dou le	a s(dou le a)	Returns the absolute value
int dou le	pmod(int a, int) pmod(dou le a, dou le)	Returns the positive value of a mod
dou le	sin(dou le a)	Returns the sine of a (a is in radians)
dou le	asin(dou le a)	Returns the arc sin of x if $-1 \leq a \leq 1$ or null otherwise
dou le	cos(dou le a)	Returns the cosine of a (a is in radians)
dou le	acos(dou le a)	Returns the arc cosine of x if $-1 \leq a \leq 1$ or null otherwise
tan(dou le a)	tan(dou le a)	Returns the tangent of a (a is in radians)
dou le	atan(dou le a)	Returns the arctangent of a
dou le	degrees(dou le a)	Converts value of a from radians to degrees
dou le	radians(dou le a)	Converts value of a from degrees to radians
int dou le	positive(int a), positive(dou le a)	Returns a
int dou le	negative(int a), negative(dou le a)	Returns -a
float	sign(dou le a)	Returns the sign of a as '1.0' or '-1.0'
dou le	e()	Returns the value of e

double

pi()

Returns the value of pi

String Functions

The following are built-in string functions supported in Hive:

Return Type	Name(signature)	Example
int	ascii(string str)	Returns the numeric value of the first character of str
string	concat(string array A, string array ...)	Returns the string or arrays resulting from concatenating the strings or arrays passed in as parameters in order. e.g. concat('foo', 'bar') results in 'foobar'. Note that this function can take an number of input strings.
array<struct<string,double>>	context_ngrams(array<array>, array, int K, int pf)	Returns the top-k contextual N-grams from a set of tokenized sentences, given a string of "context". See statisticsAndDataMining for more information.
string	concat_ws(string P, string A, string ...)	Like concat() above, but with custom separator P.
string	concat_ws(string P, array)	Like concat_ws() above, but taking an array of strings. (as of Hive 0.9.0)
int	find_in_set(string str, string	Returns the first occurrence of str in strList

	strList)	where strList is a comma-delimited string. Returns null if either argument is null. Returns 0 if the first argument contains an commas. e.g. find_in_set('a ', 'a c, ,a ,c,def') returns 3
string	format_number(num x, int d)	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	get_json_object(string json_string, string path)	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. NOT : The json path can only have the characters [0-9a-z_], i.e., no upper-case or special characters. Also, the keys *cannot start with numbers.* This is due to restrictions on Hive column names.
boolean	in_file(string str, string filename)	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 occurrence of substr in str
int	length(string A)	Returns the length of the string
int	locate(string substr, string	Returns the position of the first occurrence

	str[, int pos])	of su str in str after position pos
string	lower(string A) lcase(string A)	
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. ltrim(' foo ar ') results in 'foo ar '
arra <struct<string,dou le>>	ngrams(arra <arra >, int N, int K, int pf)	Returns the top-k N-grams from a set of tokenized sentences, such as those returned by the sentences() UDAF. See statisticsAndDataMining for more information.
string	parse_url(string url string, string partTo xtract [, string ke To xtract])	Returns the specified part from the URL. Valid values for partTo xtract include HO T, PATH, QU RY, R F, PROTOCOL, AUTHORITY, FIL , and U RINFO. e.g. parse_url('http://face ook.com/path1/p.php?k1=v1&k2=v2#Ref1', 'HO T') returns 'face ook.com'. Also a value of a particular ke in QU RY can be extracted providing the ke as the third argument, e.g. parse_url('http://face ook.com/path1/p.php?k1=v1&k2=v2#Ref1', 'QU RY', 'k1') returns 'v1'.

string	printf(string format, Object... args)	Returns the input formatted according to 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('foothe bar', 'foo(.?)(ar)', 2) returns ' ar.' Note that some care is necessary in using predefined character classes: using '\s' as the second argument will match the letter s; '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("foo bar", "oobar", "") returns 'f .' Note that some care is necessary in using predefined character classes: using '\s' as the second argument will match the letter s; '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(' foo ar ') results in ' foo ar'
array <array >	sentence(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. The 'lang' and 'locale' are optional arguments. e.g. sentence('Hello there! How are you?') returns (("Hello", "there"), ("How", "are", "you"))
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>	str_to_map(text[, delimiter1, delimiter2])	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	su str(string <array> A, int start) su string(string <array> A, int start)	Returns the sub string or slice of the array A starting from start position till the end of string A e.g. su str('foo bar', 4) results in ' bar'

string	<code>su str(string inar A, int start, int len)</code> <code>su string(string inar A, int start, int len)</code>	Returns the su string or slice of the te arra of A starting from start position with length len e.g. <code>su str('foo ar', 4, 1)</code> results in ''
string	<code>translate(string input, string from, string to)</code>	Translates the input string replacing the characters present in the from string with the corresponding characters in the to string. This is similar to the <code>translate</code> function in Postgre QL. If any of the parameters to this UDF are NULL, the result is NULL as well (available as of Hive 0.10.0)
string	<code>trim(string A)</code>	Returns the string resulting from trimming spaces from both ends of A e.g. <code>trim(' foo ar ')</code> results in 'foo ar'
string	<code>upper(string A)</code> <code>ucase(string A)</code>	Returns the string resulting from converting all characters of A to upper case e.g. <code>upper('fOo aR')</code> results in 'FOO AR'

Collection Functions

The following built-in collection functions are supported in hive:

Return Type	Name(signature)	xample
int	<code>size(Map)</code>	Returns the number of elements in the map type

int	size(Arra)	Returns the number of elements in the array type
array	map_keys(Map)	Returns an unordered array containing the keys of the input map
array	map_values(Map)	Returns an unordered array containing the values of the input map
boolean	array_contains(Array , value)	Returns TRUE if the array contains value
array	sort_array (Array)	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)

Built-in Aggregate Functions (UDAF)

The following are built-in aggregate functions supported in Hive:

Return Type	Name(signature)	Example
integer	count(*), count(expr), count(DISTINCT expr[, expr...])	count(*) – Returns the total number of retrieved rows, including rows containing NULL values; count(expr) – Returns the number of rows for which the supplied expression is non-NULL; count(DISTINCT expr[, expr]) – Returns the number of rows for which the supplied expression(s) are unique and non-NULL.
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

dou le	min(col)	Returns the minimum of the column in the group
dou le	max(col)	Returns the maximum value of the column in the group
dou le	variance(col), var_pop(col)	Returns the variance of a numeric column in the group
dou le	var_samp(col)	Returns the unbiased sample variance of a numeric column in the group
dou le	stddev_pop(col)	Returns the standard deviation of a numeric column in the group
dou le	stddev_samp(col)	Returns the unbiased sample standard deviation of a numeric column in the group
dou le	covar_pop(col1, col2)	Returns the population covariance of a pair of numeric columns in the group
dou le	covar_samp(col1, col2)	Returns the sample covariance of a pair of a numeric columns in the group
dou le	corr(col1, col2)	Returns the Pearson coefficient of correlation of a pair of a numeric columns in the group
dou le	percentile(IGBT col, p)	Returns the exact p^{th} percentile of a column in the group (does not work with floating point types). p must be between 0 and 1. NOTE : A true percentile can only be computed for integer values. Use PERCENTILE_APPROX if our input is non-integral.
arra	percentile(IGBT col, arra (p1 [, p2]...))	Returns the exact percentiles p1, p2, ... of a column in the group (does not work with floating point types). pi must be

		between 0 and 1. NOT : A true percentile can only be computed for integer values. Use PERCENTILE_APPROX if our input is non-integral.
dou le	percentile_approx(DOUBLE col, p [,])	Returns an approximate p^{th} percentile of a numeric column (including floating point types) in the group. The parameter controls approximation accuracy at the cost of memory. Higher values yield better approximations, and the default is 10,000. When the number of distinct values in col is smaller than , this gives an exact percentile value.
arra	percentile_approx(DOUBLE col, arra (p1 [, p2]...) [,])	ame as above, but accepts and returns an array of percentile values instead of a single one.
arra	histogram_numeric(col,)	Computes a histogram of a numeric column in the group using non-uniformly spaced bins. The output is an array of size of double-valued (x,) coordinates that represent the bin centers and heights
arra	collect_set(col)	Returns a set of objects with duplicate elements eliminated

Built-in Table-Generating Functions (UDTF)

Normal user-defined functions, such as concat(), take in a single input row and output a single output row. In contrast, table-generating functions transform a single input row to multiple output rows.

Return Type	Name(signature)	xample
inline(ARRAY< STRUCT[, STRUCT]>)	xplos es an array of structs into a table (as of Hive 0.10)	

xplode

explode() takes an array as an input and outputs the elements of the array as separate rows. UDTF's can be used in the LCT expression list and as a part of LAT RAL VI W.

Conditional Functions

Return Type	Name(signature)	Example
T	if(boolean testCondition, T valueTrue, T valueFalseOrNull)	Return valueTrue when testCondition is true, returns valueFalseOrNull otherwise
T	COALESCE (T v1, T v2, ...)	Return the first v that is not NULL, or NULL if all v's are NULL
T	CASE WHEN N THENC [WHEN D] WHEN A = , returns C; when A = D, return E; else F	THENC]* [ELSE F] END
T	CASE WHEN A = TRUE, Returns ; when C = TRUE, return D; ELSE E END	WHEN D]* [ELSE E] END

Functions for

Return Type

Text Analytics

Name(signature)

Example

array <struct<string, double>>

context_ngrams(array <array>, array , int K, int pf)

Returns the top-k contextual N-grams from a set of tokenized sentences, given a string of "context". See [statisticsAndDataMining](#) for

		<p>more information.N-grams are sequences 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.</p>
arra <struct<string,dou le>>	ngrams(arra <arra >, int N, int K, int pf)	<p>Returns the top-k N-grams from a set of tokenized sentences, such as those returned by the sentences() UDAF.</p> <p>See statisticsAndDataMining for more information.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 manually stripping sentences of non-contextual content and then passing them to ngrams(), but context_ngrams() makes it much easier.</p>

