

SQL TO PIG

Cheat Sheet

We know that lots of people come to Apache Pig from a relational database background, so we compiled this handy translation from SQL concepts to their Pig equivalents.

SQL CONCEPT	SQL	PIG
SELECT	<code>SELECT column_name,column_name FROM table_name;</code>	<code>FOREACH alias GENERATE column_name, column_name;</code>
SELECT *	<code>SELECT * FROM table_name;</code>	<code>FOREACH alias GENERATE *;</code>
DISTINCT	<code>SELECT DISTINCT column_name,column_name FROM table_name;</code>	<code>DISTINCT(FOREACH alias GENERATE column_name, column_name);</code>
WHERE	<code>SELECT column_name,column_name FROM table_name WHERE column_name operator value;</code>	<code>FOREACH (FILTER alias BY column_name operator value) GENERATE column_name, column_name;</code>
AND/OR	<code>... WHERE (column_name operator value1 AND column_name operator value2) OR column_name operator value3;</code>	<code>FILTER alias BY (column_name operator value1 AND column_name operator value2) OR column_name operator value3;</code>
ORDER BY	<code>... ORDER BY column_name ASC DESC, column_name ASC DESC;</code>	<code>ORDER alias BY column_name ASC DESC, column_name ASC DESC;</code>
TOP/LIMIT	<code>SELECT TOP number column_name FROM table_name ORDER BY column_name ASC DESC;</code>	<code>TOP(number, column_index, alias);</code>
	<code>SELECT column_name FROM table_name ORDER BY column_name ASC DESC LIMIT number;</code>	<code>FOREACH (GROUP alias BY column_name) GENERATE LIMIT alias number;</code>
GROUP BY	<code>SELECT function(column_name) FROM table GROUP BY column_name;</code>	<code>FOREACH (GROUP alias BY column_name) GENERATE function(alias.column_name);</code>
LIKE	<code>... WHERE column_name LIKE pattern;</code>	<code>FILTER alias BY REGEX_EXTRACT(column_name, pattern, 1) IS NOT NULL;</code>
IN	<code>... WHERE column_name IN (value1,value2,...);</code>	<code>FILTER alias BY column_name IN (value1, value2,...);</code>

SQL CONCEPT	SQL	PIG
JOIN	<code>SELECT column_name(s) FROM table1 JOIN table2 ON table1.column_name=table2.column_name;</code>	<code>FOREACH (JOIN alias1 BY column_name, alias2 BY column_name) GENERATE column_name(s);</code>
LEFT/RIGHT/FULL OUTER JOIN	<code>SELECT column_name(s) FROM table1 LEFT RIGHT FULL OUTER JOIN table2 ON table1.column_name=table2.column_name;</code>	<code>FOREACH (JOIN alias1 BY column_name LEFT RIGHT FULL, alias2 BY column_name) GENERATE column_name(s);</code>
UNION ALL	<code>SELECT column_name(s) FROM table1 UNION ALL SELECT column_name(s) FROM table2;</code>	<code>UNION alias1, alias2;</code>
AVG	<code>SELECT AVG(column_name) FROM table_name;</code>	<code>FOREACH (GROUP alias ALL) GENERATE AVG(alias.column_name);</code>
COUNT	<code>SELECT COUNT(column_name) FROM table_name;</code>	<code>FOREACH (GROUP alias ALL) GENERATE COUNT(alias);</code>
COUNT DISTINCT	<code>SELECT COUNT(DISTINCT column_name) FROM table_name;</code>	<code>FOREACH alias { unique_column = DISTINCT column_name; GENERATE COUNT(unique_column); };</code>
MAX	<code>SELECT MAX(column_name) FROM table_name;</code>	<code>FOREACH (GROUP alias ALL) GENERATE MAX(alias.column_name);</code>
MIN	<code>SELECT MIN(column_name) FROM table_name;</code>	<code>FOREACH (GROUP alias ALL) GENERATE MIN(alias.column_name);</code>
SUM	<code>SELECT SUM(column_name) FROM table_name;</code>	<code>FOREACH (GROUP alias ALL) GENERATE SUM(alias.column_name);</code>
HAVING	<code>... HAVING aggregate_function(column_name) operator value;</code>	<code>FILTER alias BY aggregate_function(column_name) operator value;</code>
UCASE/UPPER	<code>SELECT UCASE(column_name) FROM table_name;</code>	<code>FOREACH alias GENERATE UPPER(column_name);</code>
LCASE/LOWER	<code>SELECT LCASE(column_name) FROM table_name;</code>	<code>FOREACH alias GENERATE LOWER(column_name);</code>
SUBSTRING	<code>SELECT SUBSTRING(column_name,start,length) AS some_name FROM table_name;</code>	<code>FOREACH alias GENERATE SUBSTRING(column_name, start, start+length) as some_name;</code>
LEN	<code>SELECT LEN(column_name) FROM table_name;</code>	<code>FOREACH alias GENERATE SIZE(column_name);</code>
ROUND	<code>SELECT ROUND(column_name,0) FROM table_name;</code>	<code>FOREACH alias GENERATE ROUND(column_name);</code>

For more tips, download the full Pig Cheat Sheet: bit.ly/pigcheat