

Pig

Arend Hintze

“this is so cumbersome!”

- Instead of programming everything in java
MapReduce or streaming: wouldn't it be wonderful to have a simpler interface?
- Problem: break down complex MapReduce tasks into simple commands
- pig does that!

approach

- Pig Latin -> high level language to program tasks
 - takes input and cues commands on that input to create output
- Pig compiler -> takes a script, creates jobs, runs them locally or distributed on HDFS

interacting with pig

- grunt shell -> enter commands directly (start: pig -x local OR: pig -X mapreduce)
- pig pigScript.pig -> executes a pig script
- put your commands inside your java program (import org.apache.pig.PigServer)
- OR: use the graphical web interface!

typical workflow

- Load data into an alias
 - *alias* = LOAD *filename* AS (...)
- Manipulate the alias using relational operators or functions
 - *new_alias* = pig_command(*old_alias*)
- Dump alias to shell, output, or file in a HDFS directory

Pig Lating Commmands Categories

- Read-Write from/to HDFS
- Diagnostics
- Data types
- Expression and functions
- Relational operators

Read-Write Operators

Table 10.2 Data read/write operators in Pig Latin

LOAD	<pre>alias = LOAD 'file' [USING function] [AS schema];</pre> <p>Load data from a file into a relation. Uses the PigStorage load function as default unless specified otherwise with the USING option. The data can be given a schema using the AS option.</p>
LIMIT	<pre>alias = LIMIT alias n;</pre> <p>Limit the number of tuples to n. When used right after <code>alias</code> was processed by an ORDER operator, LIMIT returns the first n tuples. Otherwise there's no guarantee which tuples are returned. The LIMIT operator defies categorization because it's certainly not a read/write operator but it's not a true relational operator either. We include it here for the practical reason that a reader looking up the DUMP operator, explained later, will remember to use the LIMIT operator right before it.</p>
DUMP	<pre>DUMP alias;</pre> <p>Display the content of a relation. Use mainly for debugging. The relation should be small enough for printing on screen. You can apply the LIMIT operation on an alias to make sure it's small enough for display.</p>
STORE	<pre>STORE alias INTO 'directory' [USING function];</pre> <p>Store data from a relation into a directory. The directory must not exist when this command is executed. Pig will create the directory and store the relation in files named <code>part-$nnnn$</code> in it. Uses the PigStorage store function as default unless specified otherwise with the USING option.</p>

example

```
data = LOAD 'movie_short.csv' using PigStorage(',') AS  
(id,name,year,rating,score);  
res = FILTER data BY (float)rating>2.0;  
DUMP res;
```


LOADING

```
alias = LOAD 'file' [USING function] [AS schema];
```

- Default: assumes data is tab-delimited
- if data has different spacers use `PigStorage('delimiter')`
- the schema can have types

```
data = LOAD 'movie_short.csv' using PigStorage(',') AS  
(id:int,name:chararray,year:int,rating:float,score:float);
```

SAVING

- DUMP just dumps the output to the command line or display
- STORE saves the content to a file
- LIMIT allows you to specify the number of tuples to be returned

saving

```
data = LOAD 'movie_short.csv' using PigStorage(',') AS  
(id,name,year,rating,score);  
res = FILTER data BY (float)rating>2.0;  
STORE res INTO 'output';
```

uses tabs as delimiter

```
data = LOAD 'movie_short.csv' using PigStorage(',') AS  
(id,name,year,rating,score);  
res = FILTER data BY (float)rating>2.0;  
STORE res INTO 'output' using PigStorage(',');
```

uses other delimiter

LIMIT

```
data = LOAD 'movie_short.csv' using PigStorage(',') AS  
(id,name,year,rating,score);  
res = FILTER data BY (float)rating>2.0;  
res_limit = LIMIT res 10;  
STORE res_limit INTO 'output' using PigStorage(',');
```

DIAGNOSTIC

DESCRIBE

```
DESCRIBE alias;  
    Display the schema of a relation.
```

EXPLAIN

```
EXPLAIN [-out path] [-brief] [-dot] [-param ...]  
[-param_file ...] alias;  
    Display the execution plan used to compute a relation. When used with a script  
name, for example, EXPLAIN myscript.pig, it will show the execution plan  
of the script.
```

ILLUSTRATE

```
ILLUSTRATE alias;  
    Display step-by-step how data is transformed, starting with a load command, to  
arrive at the resulting relation. To keep the display and processing manageable,  
only a (not completely random) sample of the input data is used to simulate the  
execution.  
    In the unfortunate case where none of Pig's initial sample will survive the  
script to generate meaningful data, Pig will "fake" some similar initial data that will  
survive to generate data for alias.
```

Atomic Data Types

Table 10.4 Atomic data types in Pig Latin

<code>int</code>	Signed 32-bit integer
<code>long</code>	Signed 64-bit integer
<code>float</code>	32-bit floating point
<code>double</code>	64-bit floating point
<code>chararray</code>	Character array (string) in Unicode UTF-8
<code>bytearray</code>	Byte array (binary object)

```
data = LOAD 'movie_short.csv' using PigStorage(',') AS  
(id:int,name:chararray,year:int,rating:float,score:float);
```

Complex Data Types

Table 10.5 Complex data types in Pig Latin

Tuple	<pre>(12.5,hello world,-2)</pre> <p>A tuple is an ordered set of fields. It's most often used as a row in a relation. It's represented by fields separated by commas, all enclosed by parentheses.</p>
Bag	<pre>{(12.5,hello world,-2),(2.87,bye world,10)}</pre> <p>A bag is an unordered collection of tuples. A relation is a special kind of bag, sometimes called an outer bag. An inner bag is a bag that is a field within some complex type.</p> <p>A bag is represented by tuples separated by commas, all enclosed by curly brackets.</p> <p>Tuples in a bag aren't required to have the same schema or even have the same number of fields. It's a good idea to do this though, unless you're handling semistructured or unstructured data.</p>
Map	<pre>[key#value]</pre> <p>A map is a set of key/value pairs. Keys must be unique and be a string (<code>chararray</code>). The value can be any type.</p>

Data Types

- A field in a tuple or a value in a map can be null or any atomic/complex type (NESTING)
 - (John , {(48, Jolly Rd, Okemos),(10, Grand,Lansing)})
- Defining a schema
 - if you leave out the field type Pig will default to byte array
 - if you leave out the name a field would be unnamed and you can reference it by it's position (\$0, \$1, \$2 ... and so on)

Loading complex data types

```
cat data;
```

```
(3,8,9) (4,5,6)
```

```
(1,4,7) (3,7,5)
```

```
(2,5,8) (9,5,8)
```

tuples are tab delimited

```
A = LOAD 'data' AS (t1:tuple(t1a:int, t1b:int,t1c:int),  
                    t2:tuple(t2a:int,t2b:int,t2c:int));
```

```
DUMP A;
```

```
((3,8,9), (4,5,6))
```

```
((1,4,7), (3,7,5))
```

```
((2,5,8), (9,5,8))
```

```
X = FOREACH A GENERATE t1.t1a,t2.$0;
```

```
DUMP X;
```

```
(3,4)
```

```
(1,3)
```

```
(2,9)
```

Expressions 1

- expressions are used in FILTER, FOREACH, GROUP and SPLIT as well as in eval functions

Constant	<code>12, 19.2, 'hello world'</code>	Constant values such as 19 or "hello world." Numeric values without decimal point are treated as <code>int</code> unless <code>l</code> or <code>L</code> is appended to the number to make it a <code>long</code> . Numeric values with a decimal point are treated as <code>double</code> unless <code>f</code> or <code>F</code> is appended to the number to make it a <code>float</code> .
Basic arithmetic	<code>+, -, *, /</code>	Plus, minus, multiply, and divide.
Sign	<code>+x, -x</code>	Negation (-) changes the sign of a number.
Cast	<code>(t)x</code>	Convert the value of <code>x</code> into type <code>t</code> .
Modulo	<code>x % y</code>	The remainder of <code>x</code> divided by <code>y</code> .

Expressions 2

Conditional	<code>(x ? y : z)</code>	Returns <code>y</code> if <code>x</code> is true, <code>z</code> otherwise. This expression must be enclosed in parentheses.
Comparison	<code>==, !=, <, >, <=, >=</code>	Equals to, not equals to, greater than, less than, etc.
Pattern matching	<code>x matches regex</code>	Regular expression matching of string <code>x</code> . Uses Java's regular expression format (under the <code>java.util.regex.Pattern</code> class) to specify <code>regex</code> .
Null	<code>x is null,</code> <code>x is not null</code>	Check if <code>x</code> is null (or not).
Boolean	<code>x and y,</code> <code>x or y</code> <code>not x</code>	Boolean and, or, not.

Built-In Functions

AVG	Calculate the average of numeric values in a single-column bag.
CONCAT	Concatenate two strings (<code>chararray</code>) or two <code>bytearrays</code> .
COUNT	Calculate the number of tuples in a bag. See <code>SIZE</code> for other data types.
DIFF	Compare two fields in a tuple. If the two fields are bags, it will return tuples that are in one bag but not the other. If the two fields are values, it will emit tuples where the values don't match.
MAX	Calculate the maximum value in a single-column bag. The column must be a numeric type or a <code>chararray</code> .
MIN	Calculate the minimum value in a single-column bag. The column must be a numeric type or a <code>chararray</code> .
SIZE	Calculate the number of elements. For a bag it counts the number of tuples. For a tuple it counts the number of elements. For a <code>chararray</code> it counts the number of characters. For a <code>bytearray</code> it counts the number of bytes. For numeric scalars it always returns 1.
SUM	Calculate the sum of numeric values in a single-column bag.
TOKENIZE	Split a string (<code>chararray</code>) into a bag of words (each word is a tuple in the bag). Word separators are space, double quote ("), comma, parentheses, and asterisk (*).
IsEmpty	Check if a bag or map is empty.

case sensitive!