



# Data Recipes

Thursday, November 14, 2013

## Linear Regression With Pig

It's a bit silly trying to motivate a discussion of linear regression. It's everywhere. Linear regression is typically the first pass step for understanding a dataset. Is there a linear relationship between my variables? Maybe. Let's try linear regression! In other, and less precise, words we just fit a damn line to it. Suffice it to say linear regression is one of those tools for data analysis that's not optional. You **must** know it.

In fact, I'm going to assume just that. That you've used linear regression in other contexts before and understand its utility. The issue at hand is how to *parallelize* linear regression. Why? Well, suppose you have billions of feature vectors in your data set, each with thousands of features (columns), and you want to use all of them because, why not? Suppose it doesn't fit on one machine. Now, there exists a project to address this specifically, [vowpal wabbit](#), which you should most certainly check out, but that I'm not going to talk about. Instead, the idea is to use Apache Pig. The reason for implementing it with Pig, rather than using an existing tool, is mostly for illustration. Linear regression with pig brings up several design and implementation details that I believe you'll face when doing almost any reasonably useful machine learning at scale. In other words, how do I wire all this shit together?

Specifically, I'll address pig macros, python drivers, and using a whole relation as a scalar. Fun stuff.

### linear regression

It's important that I do at least explain a bit of terminology so we're all together in this. So, rather than jump for the most general explanation immediately (why do people do that?) let's talk about something real. Suppose you've measured the current running through a circuit while slowly decreasing the resistance. You would expect the current to increase linearly as you decrease the current ([ohms law](#)). In other words,

$$I = \frac{V}{R}$$

To verify Ohm's Law (not a bad idea, I mean, maybe you're living in a dream world where physics is different and you want to know for certain...) you'd record the current  $I$  and the resistance  $R$  at every measurement while holding the voltage  $V$  constant. You'd then fit a line to the data or, more specifically, to  $\frac{1}{R}$ , and, if all went well, find that the slope of said line was equal to the voltage.

Terminology: the current would be called the *response* or target variable. All the responses together form a vector  $y$  called the *response vector*. The resistance would be called the *feature* or observation. And, if you recorded more than just the resistance, say, the temperature, then for every response you'd have a collection of features or a *feature vector*. All the feature vectors together form a matrix  $X$ . The goal of linear regression is to find the best set of weights  $w$  that, when used to form a linear combination of the features, creates a vector that is as close as possible to the response vector. So the problem can be phrased as an optimization problem. Here the function we'll minimize is the mean squared error (the square of the distance between the weighted features and the response). Mathematically the squared error, for one feature vector  $x_i$  of length  $M$ , can be written as:

$$error^2 = (y_i - \sum_{j=1}^M w_j x_{i,j})^2$$

where  $x_{i,0} = 1$  by definition.

So the mean squared error (mse), when we've got  $N$  measurements (feature vectors), is then:

$$mse(w) = \frac{1}{N} \sum_{i=1}^N (y_i - \sum_{j=1}^M w_j x_{i,j})^2$$

So now the question is, exactly how are we going to minimize the mse by varying the weights? Well, it turns out there's the method called *gradient descent*. That is, the mse decreases fastest if we start with a given set of weights and travel in the direction of the negative gradient of the mse for those weights. In other words:

$$w_{new} = w - \alpha \nabla mse(w)$$

Where  $\alpha$  is the magnitude of the step size. What this gives us is a way to update the weights until  $w_{new}$  doesn't really change much. Once the weights converge we're done.

### algorithm

Alright, now that we've got a rule for updating weights, we can write down the algorithm.

#### Followers

#### Blog Archive

- ▼ 2013 (3)
  - ▼ November (1)
    - Linear Regression With Pig
  - September (1)
  - August (1)
- 2012 (2)
- 2011 (21)

#### About Me



thedatachef

@thedatachef

[View my complete](#)

Loading [MathJax]/jax/output/HTML-CSS/jax.js

1. Initialize the weights, one per feature, randomly
2. Update the weights by subtracting the gradient of the mse
3. Repeat 2 until converged

## implementation

Ok great. Let's get the pieces together.

### pig

Pig is going to do most of the real work. There's two main steps involved. The first step, and one that varies strongly from domain to domain, problem to problem, is loading your data and transforming it into something that a generic algorithm can handle. The second, and more interesting, is the implementation of gradient descent of the mse itself.

Here's the implementation of the gradient descent portion. I'll go over each relation in detail.

```

1  /*
2   * Run one iteration of gradient descent for the given weights and features
3   * with step size alpha. Returns the updated weights.
4   *
5   * features - Relation with the following schema:
6   * {response:double, vector:tuple(f0:double,f1:double,...,fN:double)}
7   *
8   * w - Relation with **exactly one tuple** with the following schema:
9   * {weights:tuple(w0:double,w1:double,...,wN:double)}
10  *
11  * alpha - Step size. Schema:
12  *   alpha:double
13  */
14  define sounder_lm_gradient_descent(features, w, alpha) returns new_weights {
15
16      --
17      -- Use a scalar cast to zip weights with their appropriate features
18      --
19      weighted = foreach $features {
20          zipped = Zip($w.weights, vector);
21          generate
22              response as response,
23              vector   as vector,
24              zipped   as zipped:bag{t:tuple(weight:double, feature:double, dimension:int)};
25      };
26
27
28      --
29      -- Compute dot product of weights with feature vectors. First part of
30      -- step adjustment.
31      --
32      dot_prod = foreach weighted {
33          dots      = foreach zipped generate weight*feature as product;
34          dot_product = SUM(dots.product);
35          diff       = (dot_product-response);
36
37          generate
38              flatten(zipped) as (weight, feature, dimension), diff as diff;
39      };
40
41
42      scaled = foreach dot_prod generate dimension, weight, feature*diff as feature_diff;
43
44      --
45      -- Compute step diff along each dimension. Uses combiners
46      --
47      steps = foreach (group scaled by (dimension,weight)) {
48          factor      = ($alpha/(double)COUNT(scaled));
49          weight_step = factor*SUM(scaled.feature_diff);
50          new_weight  = group.weight - weight_step;
51          generate
52              group.dimension as dimension,
53              new_weight      as weight;
54      };
55
56      --
57      -- A group all is acceptable here since the previous step reduces the
58      -- size down to the number of features.
59      --
60      $new_weights = foreach (group steps all) {
61          in_order = order steps by dimension asc;
62          as_tuple = BagToTuple(in_order.weight);
63          generate

```

```

64         as_tuple;
65     };
66 };

```

linear\_model.pig hosted with ❤ by GitHub

[view raw](#)

Go ahead and save that in a directory called 'macros'. Since gradient descent of the mean squared error for the purposes of creating a linear model is a generic problem, it makes sense to implement it as a pig macro.

In lines 19-25 we're attaching the weights to every feature vector. The Zip UDF, which can be found on [github](#), receives the weights as a tuple and the feature vector as a tuple. The output is a bag of new tuples which contains (weight,feature,dimension). Think of Zip like a zipper where it matches the weights to their corresponding features. Importantly, the dimension (index in the input tuples) is returned as well.

Something to notice about this first bit is the scalar cast. Zip receives the *entire* weights relation as a single object. Since the weights relation is only a single tuple anyway, this is great. This is a good thing. It prevents us from doing something silly like a replicated join on a constant (which works but clutters the logic) or, worse, a cross.

Next, in lines 32-39, we're computing a portion of the gradient of the mse. The reason for the Zip udf in the first step was so the nested projection and sum to compute the dot product of the weights with the features works out cleanly.

Then, on line 42, the full gradient of the mse is computed by multiplying each feature by the error. It might not seem obvious when written like that, but the partial derivatives of the mse with respect to each weight make it work out like this. How nice.

Lines 47-54 is where the action happens. By action I mean we'll actually trigger a reduce job since everything up to this point has been map only. This is where the partial derivative bits come together. That is, we're grouping by dimension and weight (which, it turns out is precisely the thing we're differentiating with respect to) and summing each feature vector's contribution to the partial derivative along that dimension. The factor bit, on line 48, is the multiplier for the gradient. It includes the normalization term (we normalize by the number of features since we're differentiating the *mean* squared error) and the step size. The result of this step is a new weight for each dimension.

An important thing to note here is that there are only a number of partitions equivalent to the number of features. What this means is that each reduce task would get a potentially *very large* amount of data. Fortunately for us COUNT and SUM are both algebraic and so Pig will use combiners, hurray!, drastically reducing the amount of data sent to each reduce task.

Finally, on lines 60-65, we reconstruct a new tuple with the new weights and return it. The schema of this tuple should be the same as the input weight tuple.

### gory details

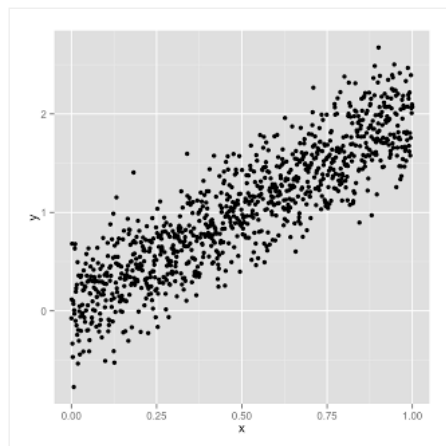
So now that we have a way to do it, we should do it. Right, I mean if we can we should. Isn't that how technology works...

I'll going to go ahead and do a completely contrived example. The reason is so that I can visualize the results.

I've created some data, called data.tsv, which satisfies the following:

$$y = 0.3r(x) + 2x$$

where  $r(x)$  is a noise term. And here's the plot:



So we have two feature columns,  $(1.0, x)$ , that we're trying to find the weights for. Since we've cooked this example (we already know the relationship between  $x$  and  $y$ ) we expect the weights for those columns to be  $(0.0, 2.0)$  if all goes well.

Now that we've got some data to work with, we'll need to write a bit more pig to load the data up and run it through our gradient descent macro.

```

1 %default STEPSIZE 0.1
2

```

```

3  import 'macros/linear_model.pig';
4
5  data = load '$data' as (x:double, y:double);
6
7  --
8  -- The weights are given a schema file by the python driver script
9  -- to avoid having to manually specify and arbitrarily large schema
10 -- with identical column schemas
11 --
12 weights = load '$input_weights' using PigStorage('\t', '-schema');
13
14
15 features = foreach data generate y as response, TOTUPLE(1.0,x) as vector;
16 weights = foreach weights generate TOTUPLE(w0,w1) as weights;
17
18 new_weights = sounder_lm_gradient_descent(features, weights, $STEPSIZE);
19 new_weights = foreach new_weights generate flatten($0);
20
21 --
22 -- Store weights in the same way they were came in. The
23 -- python driver copies the input weights schema to the output
24 -- dir (since it's the same). The reason for this is to avoid
25 -- having to manually specify an as clause, in the project
26 -- expression above, with over 100 fields.
27 --
28 store new_weights into '$output_weights';

```

fit\_line.pig hosted with ❤ by GitHub

[view raw](#)

There's really not much exciting going on here. We're loading the features and weights and rearranging them to satisfy the schema that the gradient descent macro expects. The other details are related to the driver script.

The driver script is next. Since our algorithm is iterative and pig itself has no support for iteration, we're going to embed the pig script into a python (jython) program. Here's what that looks like:

```

1  #!/usr/bin/env python
2
3  import os
4  import sys
5  import shutil
6  import random
7  import tempfile
8  import Queue
9
10 from org.apache.pig.scripting import Pig
11 from org.codehaus.jackson.map import ObjectMapper
12
13 EPS = 10e-6 # maximum distance between consecutive weights for convergence
14
15 pig_script = sys.argv[1] # pig script to run iteratively
16 data_dir = sys.argv[2] # directory where intermediate weights will be written
17 features = sys.argv[3] # location, inside data_dir, where the data to fit exists
18 num_features = sys.argv[4] # number of features
19
20 #
21 # Cleanup data dir
22 #
23 cmd = "rmr %s/weight-*" % data_dir
24 Pig.fs(cmd)
25
26 #
27 # Initialize weights
28 #
29 w0_fields = []
30 weights = []
31 for i in xrange(int(num_features)):
32     weights.append(str(random.random()))
33     w0_fields.append({"name": "w%s" % i, "type": 25, "schema": None}) # See Pig's DataType.java
34
35 path = tempfile.mkdtemp()
36 w0 = open("%s/part-r-00000" % path, 'w')
37 w0.write("\t".join(weights)+"\n")
38 w0.close()
39
40 #
41 # Create schema for weights, place under weight-0 dir
42 #
43 w0_schema = {"fields": w0_fields, "version": 0, "sortKeys": [], "sortKeyOrders": []}

```

```

44 w0_schema_file = open("%s/.pig_schema" % path, 'w')
45 ObjectMapper().writeValue(w0_schema_file, w0_schema);
46 w0_schema_file.close()
47
48 #
49 # Copy initial weights to fs
50 #
51 copyFromLocal = "copyFromLocal %s %s/%s" % (path, data_dir, "weight-0")
52 Pig.fs(copyFromLocal)
53
54 #
55 # Iterate until converged
56 #
57 features = "%s/%s" % (data_dir, features)
58 script = Pig.compileFromFile(pig_script)
59 weight_queue = Queue.Queue(25) # for moving average
60 avg_weight = [0.0 for i in xrange(int(num_features))]
61 converged = False
62 prev = 0
63 weight_dir = tempfile.mkdtemp()
64
65 while not converged:
66     input_weights = "%s/weight-%s" % (data_dir, prev)
67     output_weights = "%s/weight-%s" % (data_dir, prev+1)
68
69     bound = script.bind({'input_weights':input_weights, 'output_weights':output_weights, 'data':feature
70     bound.runSingle()
71
72 #
73 # Copy schema for weights to each output
74 #
75 copyOutputSchema = "cp %s/.pig_schema %s/.pig_schema" % (input_weights, output_weights)
76 Pig.fs(copyOutputSchema)
77
78 #
79 # The first few iterations the weights bounce all over the place
80 #
81 if (prev > 1):
82     copyToLocalPrev = "copyToLocal %s/part-r-00000 %s/weight-%s" % (input_weights, weight_dir, pr
83     copyToLocalNext = "copyToLocal %s/part-r-00000 %s/weight-%s" % (output_weights, weight_dir, p
84
85     Pig.fs(copyToLocalPrev)
86     Pig.fs(copyToLocalNext)
87
88     localPrev = "%s/weight-%s" % (weight_dir, prev)
89     localNext = "%s/weight-%s" % (weight_dir, prev+1)
90
91     x1 = open(localPrev, 'r').readlines()[0]
92     x2 = open(localNext, 'r').readlines()[0]
93
94     x1 = [float(x.strip()) for x in x1.split("\t")]
95     x2 = [float(x.strip()) for x in x2.split("\t")]
96
97     weight_queue.put(x1)
98
99     avg_weight = [x[1] + (x[0] - x[1])/(prev-1.0) for x in zip(x1, avg_weight)]
100
101 #
102 # Make sure to collect enough weights into the average before
103 # checking for convergence
104 #
105 if (prev > 25):
106     first_weight = weight_queue.get()
107     avg_weight = [(x[0] - x[1])/25.0 + x[2]/25.0 for x in zip(avg_weight, first_weight, x1)]
108
109 #
110 # Compute distance from weight centroid to new weight
111 #
112 d = sum([(pair[0] - pair[1])**2 for pair in zip(x2, avg_weight)])
113
114 converged = (d < EPS)
115
116 os.remove(localPrev)
117 os.remove(localNext)
118
119
120
121

```



```

122     prev += 1
123
124     #
125     # Cleanup
126     #
127     shutil.rmtree(path)
128     shutil.rmtree(weight_dir)

```

driver.py hosted with ❤ by GitHub

[view raw](#)

There's a lot going on here (it's cluttered looking because Pig doesn't allow class or function definitions in driver scripts) that's not that interesting and pretty easy to understand so I'll just go over the high points:

- We initialize the weights randomly and write a .pig\_schema file as well. The reason for writing the schema file is so that it's unnecessary to write the redundant schema for weights in the pig script itself.
- We want the driver to be agnostic to whether we're running in local mode or mapreduce mode. Thus we copy the weights to the filesystem using the *copyFromLocal* fs command. In local mode this just puts the initial weights on the local fs whereas in mapreduce mode this'll place them on the hdfs.
- Then we iterate until the convergence criteria is met. Each iteration details like copying the schema, and pulling down the weights to compute distances is done.
- A moving average is maintained over the past 25 weights. Iteration stops when the new weight is less than EPS away from the average.

Aside from that there's an interesting -bug?- that comes up as a result of this. Notice on line 70 how the new variables are bound each iteration? Well the underlying PigContext object just keeps adding these new variables instead of overwriting them. What this means is that, after a couple thousand iterations, depending on your PIG\_HEAPSIZE env variable, the driver script will crash from an out of memory error. Yikes.

## run it

Finally we get to run it. That part's easy:

```

$: export PIG_HEAPSIZE=4000
$: pig driver.py fit_line.pig data data.tsv 2

```

That is, we use the pig command to launch the driver program. The arguments to the driver script itself follow. We're running the fit\_line.pig script where the data dir (where intermediate weights will go) exists under 'data' and the input data, data.tsv, should exist there as well. The '2' indicates we've got two weights (w0, and w1). The pig heapsize is to deal with the bug mentioned in the previous section.

On my laptop, in local mode, the convergence criteria was met after 1527 iterations.

## results

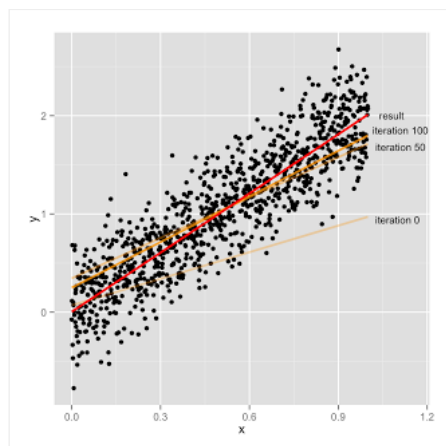
After 1527 iterations the weights ended up as (0.0040666759502969215, 2.0068029003414014) which is exactly what we'd expect. In other words:

$$y = 0.0041 + 2.0068x$$

which is the 'best fit' line to our original:

$$y = 0.3r(x) + 2x$$

And here's an illustration of what happened.



Looks pretty reasonable to me. Hurray. Now go fit some real data.

Posted by thedatachef at 10:02 AM



Labels: [apache hadoop](#), [apache pig](#), [hadoop](#), [hadoop regression](#), [linear regression](#), [pig](#), [pig examples](#), [pig python](#), [pig](#)

## 2 comments:



**rjurney** November 14, 2013 at 7:28 PM

Awesomiddityawesome!

[Reply](#)



**han yufei** June 18, 2014 at 12:52 AM

I've tried it on the cluster but it doesn't work :( When it reaches 'zipped = Zip(\$w.weights, vector);', I get the runtime exception error 'ERROR 0: Scalar has more than one row in the output. That is very confused. weight.weights is just one single tuple, like vector, why can't we use it as input argument of udf ?

,

[Reply](#)

[Home](#)

[Older Post](#)

Subscribe to: [Post Comments \(Atom\)](#)