

# BigML Assignment 5: Logistic Regression Using Stochastic Gradient Descent

Due: Wednesday, Mar. 26, 2014 13:29 EST via Autolab

Late submission with 50% credit: Friday, Mar. 28, 2014 13:29  
EST via Autolab

## Policy on Collaboration among Students

These policies are the same as were used in Dr. Rosenfeld's previous version of 10601 from 2013. The purpose of student collaboration is to facilitate learning, not to circumvent it. Studying the material in groups is strongly encouraged. It is also allowed to seek help from other students in understanding the material needed to solve a particular homework problem, provided no written notes are shared, or are taken at that time, and provided learning is facilitated, not circumvented. The actual solution must be done by each student alone, and the student should be ready to reproduce their solution upon request. The presence or absence of any form of help or collaboration, whether given or received, must be explicitly stated and disclosed in full by all involved, on the first page of their assignment. Specifically, each assignment solution must start by answering the following questions in the report:

- Did you receive any help whatsoever from anyone in solving this assignment? Yes / No. If you answered 'yes', give full details: \_\_\_\_\_ (e.g. "Jane explained to me what is asked in Question 3.4")
- Did you give any help whatsoever to anyone in solving this assignment? Yes / No. If you answered 'yes', give full details: \_\_\_\_\_ (e.g. "I pointed Joe to section 2.3 to help him with Question 2").

Collaboration without full disclosure will be handled severely, in compliance with CMU’s Policy on Cheating and Plagiarism. As a related point, some of the homework assignments used in this class may have been used in prior versions of this class, or in classes at other institutions. Avoiding the use of heavily tested assignments will detract from the main purpose of these assignments, which is to reinforce the material and stimulate thinking. Because some of these assignments may have been used before, solutions to them may be (or may have been) available online, or from other people. It is explicitly forbidden to use any such sources, or to consult people who have solved these problems before. You must solve the homework assignments completely on your own. I will mostly rely on your wisdom and honor to follow this rule, but if a violation is detected it will be dealt with harshly. Collaboration with other students who are currently taking the class is allowed, but only under the conditions stated below.

## 1 Important Note

**As usual, you are expected to use Java for this assignment. It could take hours to run your experiments. Start early.**

Please post clarification questions to the Piazza, and the instructors can be reached at the following email address: *10605-Instructors@cs.cmu.edu*.

## 2 Background: SGD for Logistic Regression

One fairly simple way (and extremely scalable way) to implement logistic regression is stochastic gradient descent.

In the lecture we followed Charles Elkan’s notes<sup>1</sup>, which are for a binary classification task. We estimate the probability  $p$  that an example  $\mathbf{x} = \langle x_1, \dots, x_d \rangle$  is positive in the log-odds form:

$$\log \frac{p}{1-p} = \alpha + \sum_{j=1 \dots d} \beta_j x_j \quad (1)$$

If we assume there is a “bias feature”  $x_0$  that is true for every example, then

---

<sup>1</sup><http://cseweb.ucsd.edu/~elkan/250B/logreg.pdf>

you can simplify and drop the  $\alpha$ , leaving just the  $\beta_j$ 's to estimate. Therefore

$$p = \frac{\exp(\beta^T \mathbf{x})}{1 + \exp(\beta^T \mathbf{x})}. \quad (2)$$

It's convenient to consider examples of the form  $\mathbf{x}, y$  where  $y = 0$  or  $y = 1$ . The log of the conditional likelihood for example will be  $LCL(\mathbf{x}, y) = \log p$  if  $y = 1$  and  $LCL(\mathbf{x}, y) = \log(1 - p)$  if  $y = 0$ , where  $p$  is computed as in Eq. 1. With a little calculus you can show that for a positive example,

$$\frac{\partial}{\partial \beta_j} LCL(\mathbf{x}, y) = \frac{1}{p} \frac{\partial}{\partial \beta_j} p$$

and for a negative example,

$$\frac{\partial}{\partial \beta_j} LCL(\mathbf{x}, y) = \frac{1}{1 - p} \left( -\frac{\partial}{\partial \beta_j} p \right)$$

and that

$$\frac{\partial}{\partial \beta_j} p = p(1 - p)x_j$$

and putting this together we get that if  $y = 1$

$$\frac{\partial}{\partial \beta_j} LCL(\mathbf{x}, y) = (1 - p)x_j$$

and if  $y = 0$  then

$$\frac{\partial}{\partial \beta_j} LCL(\mathbf{x}, y) = -px_j$$

so in either case

$$\frac{\partial}{\partial \beta_j} LCL(\mathbf{x}, y) = (y - p)x_j \quad (3)$$

So an update to the  $\beta$ 's that would improve most LCL would be along the gradient—i.e., for some small step size  $\lambda$ , let

$$\beta_j = \beta_j + \lambda(y - p)x_i$$

Notice that if  $x_i = 0$  then  $\beta_j$  is unchanged.

So this leads to this algorithm, which is very fast (assuming you have enough memory to hash all the parameter values).

1. Initialize a hashtable  $B$
2. For  $t = 1, \dots, T$ 
  - For each example  $\mathbf{x}_i, y_i$ :
    - For each non-zero feature of  $\mathbf{x}_i$  with index  $j$  and value  $x_j$ :
      - \* If  $j$  is not in  $B$ , set  $B[j] = 0$ .
      - \* Set  $B[j] = B[j] + \lambda(y - p)x_i$
3. Output the parameters  $\beta_1, \dots, \beta_d$ .

The time to run this is  $O(nT)$ , where  $n$  is the total number of non-zero features for each example and  $T$  is the number of iterations.

### 3 Efficient regularized SGD

Logistic regression tends to overfit when there are many rare features. One fix is to penalize large values of  $\beta$ , by optimizing, instead of  $LCL$ , some function such as  $LCL - \mu \sum_{j=1}^d \beta_j^2$ . Here  $\mu$  controls how much weight to give to the penalty term. The update for  $\beta_j$  becomes

$$\beta_j = \beta_j + \lambda((y - p)x_i - 2\mu\beta_j)$$

or equivalently

$$\beta_j = \beta_j + \lambda(y - p)x_i - \lambda 2\mu\beta_j$$

Experimentally this greatly improves overfitting - but unfortunately, this makes the computation much more expensive, because now *every*  $\beta_j$  needs to be updated, not only the ones that are non-zero.

The trick to making this efficient is to break the update into two parts. One is the usual update of adding  $\lambda(y - p)x_i$ . Let's call this the "LCL" part of the update. The second is the "regularization part" of the update, which is to replace  $\beta$  by

$$\beta_j = \beta_j - \lambda 2\mu\beta_j = \beta_j \cdot (1 - 2\lambda\mu)$$

So we could perform our update of  $\beta_j$  as follows:

- Set  $\beta_j = \beta_j \cdot (1 - 2\lambda\mu)$

- If  $x_j \neq 0$ , set  $\beta_j = \beta_j + \lambda(y - p)x_i$

Following this up, we note that we can perform  $m$  successive “regularization” updates by letting  $B_j = B_j \cdot (1 - 2\lambda\mu)^m$ . The basic idea of the new algorithm is to not perform regularization updates for zero-valued  $x_j$ ’s, but instead to simply keep track of how many such updates would need to be performed to update  $\beta_j$ , and perform them only when we would normally perform “LCL” updates (or when we output the parameters at the end of the day).

Here’s the final algorithm (for more detail, see “Lazy sparse stochastic gradient descent for regularized multinomial logistic regression”, Bob Carpenter<sup>2</sup>)

1. Let  $k = 0$ , and let  $A$  and  $B$  be empty hashables.  $A$  will record the value of  $k$  last time  $B[j]$  was updated.
2. For  $t = 1, \dots, T$ 
  - For each example  $\mathbf{x}_i, y_i$ :
  - Let  $k = k + 1$ 
    - For each non-zero feature of  $\mathbf{x}_i$  with index  $j$  and value  $x_j$ :
      - \* If  $j$  is not in  $B$ , set  $B[j] = 0$ .
      - \* If  $j$  is not in  $A$ , set  $A[j] = 0$ .
      - \* Simulate the “regularization” updates that would have been performed for the  $k - A[j]$  examples since the last time a non-zero  $x_j$  was encountered by setting
$$B[j] = B[j] \cdot (1 - 2\lambda\mu)^{k - A[j]}$$
      - \* Set  $B[j] = B[j] + \lambda(y - p)x_i$
      - \* Set  $A[j] = k$
3. For each parameter  $\beta_1, \dots, \beta_d$ , set

$$B[j] = B[j] \cdot (1 - 2\lambda\mu)^{k - A[j]}$$

4. Output the parameters  $\beta_1, \dots, \beta_d$ .

---

<sup>2</sup><http://lingpipe.files.wordpress.com/2008/04/lazysgdregression.pdf>

The learning rate  $\lambda$  is often decreased over time. On the  $t$ -th sweep through the data, set  $\lambda = \frac{\eta}{t^2}$ . I used  $\eta = 0.5$ . (Sometimes  $\lambda$  is also scaled by  $1/n_e$ , where  $n_e$  is the number of examples.) I also used a value of  $2\mu = 0.1$ .

When running stochastic gradient descent, it is usual to randomize the order of examples, and scale the feature values so that they are comparable (if they are not already binary). However, randomization is not trivial to do for a large dataset. I recommend implementing SGD as a process that streams once through a data stream, with the number of examples  $n_e$  being passed in separately as a command-line argument so that the algorithm is aware of what the current value of  $t$  is. Then write a separate module that will input a file of examples and then stream the individual examples out in approximately random order. **Hint:** to randomly sort a file in linux you can do

```
cat -n text_file | sort -R | cut -f2-
```

## 4 Task

For this assignment, we will be classifying Wikipedia articles by the languages in which they are also available. The data appears at `/afs/cs.cmu.edu/project/bigML/dbpedia/`

The data format is the same as for Assignment 1. There is one instance per line. The first token of each line is the (comma separated) list of class names, then a tab, then the data.

To reduce the experimental load we will only use the abstract data sets for this assignment. Again, they are in increasing size so that you can debug your code on smaller data. The files that start with *abstract* include Wikipedia text.

```
abstract.small.test
abstract.small.train
abstract.test
abstract.train
abstract.tiny.test
abstract.tiny.train
```

In contrast to the previous assignments, we are going to treat the multilabel classification problem as 14 independent binary classification tasks.

This means that we are going to train 14 binary classifiers<sup>3</sup>. Ideally, to obtain perfect performance (which is often difficult in practice), you will need to correctly predict all labels associated with the target document.

To reduce the experimental load fix the number of training iterations (scans of data sets) to 20 for all data sets. To fit our model to the memory of a desktop we will use the hash trick discussed in class: map every word to a features id in the range  $0 - N$ , where  $N$  is the dictionary size.

**Hint:** to convert a string to an id between 0 and  $N$  you can do something like

```
int id = word.hashCode() % N;
if (id<0) id+= N;
```

## 5 Deliverables

Submit your implementations via AutoLab. You should implement the algorithm by yourself instead of using any existing machine learning toolkit. You should upload your code (including all your function files) along with a **report**, which should solve the following questions:

1. Show values of overall likelihood function for each iteration when training with the small data set having dictionary size 10000 and  $\mu = 0.1$ . The objective function is defined as the sum of all 14 classes  $\sum_i \sum_c LCL_c(\mathbf{x}^i, y^{c,i})$ . Here  $c$  is the label id and  $i$  is the document id (5 points).

**Hint:** in order to prevent overflow when calculating  $p$  as defined in equation (2) you can use a special version of sigmoid function as the following:

```
static double overflow=20;
protected double sigmoid(double score) {
    if (score > overflow) score =overflow;
    else if (score < -overflow) score = -overflow;
    double exp = Math.exp(score);
    return exp / (1 + exp);
}
```

---

<sup>3</sup>one for each class nl,el,ru,sl,pl,ca,fr,tr,hu,de,hr,es,ga,pt.

2. Show accuracy curves for the full data set with varying regularization parameter  $\mu = 0, 1e-6, 1e-5, 1e-4, 1e-3, 0.01, 0.1, 0.2, 0.3, 0.5, 1$  and fixed dictionary size  $1e5$ , and discuss (5 points).
3. Show accuracy curves for the full data set with varying dictionary sizes  $D = 10, 100, 1e3, 1e4, 1e5, 1e6$  and the best  $\mu$  values you found in the previous step, and discuss (5 points).
4. Compare SGD logistic regression with naive Bayes classifier and discuss. (You will need to retest the naive Bayes classifiers on all three datasets with the new evaluation metric in this assignment) (5 points).
5. Answer the questions in the collaboration policy on page 1.

## Autolab Implementation Details

Your logistic regression program LR.java should be able to run without the out-of-memory issue, using the follow example command:

```
for((i=1;i<=20;i++));
do shuf trainData.txt;
done | java -Xmx128m LR 10000 0.5 0.1 20 1000 testData.txt
```

Here, the first argument 10000 is the vocabulary size. Second argument 0.5 is the initial value of the learning rate. The third argument 0.1 is the regularization coefficient. The fourth argument 20 is the max iteration (# of passes through data). The fifth argument is the size of the training dataset (which allows you to determine the starting point of a new pass). The command produces output in the following format (**one-line-per-test-example**):

```
nl<tab>p_nl,el<tab>p_el,ru<tab>p_ru,sl<tab>p_sl,pl<tab>p_pl,...
ca<tab>p_ca,fr<tab>p_fr,tr<tab>p_tr,hu<tab>p_hu,de<tab>p_de,...
hr<tab>p_hr,es<tab>p_es,ga<tab>p_ga,pt<tab>p_pt
```

Here, binary classifiers and their prediction scores are separated by comma. For each of the classifier, for example, in the first entry, nl refers to the name of the classifier, and p\_nl is the posterior probability of the label given the observation. If the posterior is larger than or equal to 0.5, then it means your classifier predicts a positive output for this target label in this document. A tab is used to separate the label name and its posterior.



You should tar the following items into **hw5.tar** and submit to the homework 5 assignment via Autolab:

- LR.java
- and all other auxiliary functions you have written
- report.pdf

Tar the files directly using “tar -cvf hw5.tar \*.java report.pdf”. Do **NOT** put the above files in a folder and then tar the folder. You do not need to upload the saved temporary files. Please make sure your code is working fine on linux.andrew.cmu.edu machines before you submit.

## 6 Submission

You must submit your homework through Autolab via the “Homework5: Memory Efficient Stochastic Gradient Descent” link. In this homework, we provide an additional tool called “Homework5-validation”:

- Homework5-validation: You will be notified by Autolab if you can successfully finish your job on the Autolab virtual machines. Note that this is not the place you should debug or develop your **algorithm**. All development should be done on linux.andrew.cmu.edu machines. This is basically a Autolab debug mode. There will be **NO** feedback on your **performance** in this mode. You have unlimited amount of submissions here. To avoid Autolab queues on the submission day, the validation link will be closed 24 hours prior to the official deadline. If you have received a score of 1000, this means that you code has passed the validation.
- Homework5: Memory Efficient Stochastic Gradient Descent: This is where you should submit your validated final submission. You have a total of **10 possible submissions**. Your performance will be evaluated, and feedback will be provided immediately.

## 7 Grading

You will be graded based on the memory usage (30 points) and runtime (30 points) for your code, and your final test performance (20 points). Note that

each pass in SGD is randomized, and there will be sampling variance when measuring your memory and runtime information. You may contact us to use your highest Autolab submission grade after homework is due. The report will be graded manually (20 points).