

CSE 6242 / CX 4242: Data and Visual Analytics | Georgia Tech | Fall 2016
Homework 4 : Scalable PageRank via Virtual Memory (MMap), Random Forest, Weka

Due: Sunday, December 4, 2016, 11:55 PM EST

Prepared by Nilaksh Das, Pradeep Vairamani, Vishakha Singh, Yanwei Zhang,
Bhanu Verma, Meghna Natraj, Polo Chau

Submission Instructions and Important Notes:

It is important that you read the following instructions carefully and also those about the deliverables at the end of each question or **you may lose points**.

- ❑ Submit a single zipped file, called “HW4-{YOUR_LAST_NAME}-{YOUR_FIRST_NAME}.zip”, containing all the deliverables including source code/scripts, data files, and readme. Example: ‘HW4-Doe-John.zip’ if your name is John Doe. Only .zip is allowed (no .rar, etc.)
- ❑ You may collaborate with other students on this assignment, but you must write your own code and give the explanations in your own words, and also mention the collaborators’ names on T-Square’s submission page. All GT students must observe [the honor code](#). **Suspected plagiarism and academic misconduct will be reported to and directly handled** by the [Office of Student Integrity \(OSI\)](#). Here are some examples similar to Prof. Jacob Eisenstein’s [NLP course page](#) (grading policy):
 - ❑ **OK:** discuss concepts (e.g., how cross-validation works) and strategies (e.g., use hashmap instead of array)
 - ❑ **Not OK:** several students work on one master copy together (e.g., by dividing it up), sharing solutions, or using solution from previous years or from the web.
- ❑ If you use any “*slip days*”, you must write down the number of days used in the T-square submission page. For example, “Slip days used: 1”. Each slip day equals 24 hours. E.g., if a submission is late for 30 hours, that counts as 2 slip days.
- ❑ At the end of this assignment, we have specified a folder structure about how to organize your files in a single zipped file. **5 points will be deducted for not following this strictly.**
- ❑ We will use auto-grading scripts to grade some of your deliverables (there are hundreds of students), so it is extremely important that you strictly follow our requirements. **Marks may be deducted if our grading scripts cannot execute on your deliverables.**
- ❑ Wherever you are asked to write down an explanation for the task you perform, **stay within the word limit** or you may lose points.
- ❑ In your final zip file, please **do not include any intermediate files** you may have generated to work on the task, unless your script is absolutely dependent on it to get the final result (which it ideally should not be).
- ❑ After all slip days are used up, **5% deduction for every 24 hours of delay**. (e.g., 5 points for a 100-point homework)
- ❑ **We will not consider late submission of any missing parts** of a homework assignment or project deliverable. To make sure you have submitted everything, download your submitted files to double check.

Task 0: Download HW skeleton

Download the HW skeleton from [this link](#).

Task 1: Scalable single-PC PageRank on 70M edge graph (40 points)

In this task, you will learn how to use your computer's [virtual memory](#) to implement the PageRank algorithm that will scale to graph datasets with [as many as billions of edges](#) using a single computer (e.g., your laptop). As discussed in class, a predominant way to work with larger datasets has been to use computer clusters (e.g., Spark, Hadoop) which may involve steep learning curves, may be costly (e.g., pay for hardware and personnel), and importantly may be “overkill” for smaller datasets (e.g., a few tens or hundreds of GBs). The virtual memory based approach offers an attractive, simple solution to allow practitioners and researchers to more easily work with such data (visit the [NSF-funded MMap project's homepage](#) to learn more about the research).

The main idea is to put the dataset in your computer's (unlimited) virtual memory, as the dataset is often too big to fit in RAM. When running algorithms on that dataset (e.g., PageRank), the operating system will automatically decide when to load the necessary data (subset of whole dataset) into RAM.

The technical approach to put data into your machine's virtual memory space is called “memory mapping”, which allows the dataset to be treated as if it is an in-memory dataset. That is, in your (PageRank) program, you do not need to know whether the data that you need is stored on the hard disk, or kept in RAM. Note that memory-mapping a file [does NOT cause the whole file to be read into memory](#). Instead, data is loaded and kept in memory only when needed (determined by strategies like [least recently used](#) paging and [anticipatory](#) paging).

You will use the Python modules [mmap](#) and [struct](#) to map a large graph dataset into your computer's virtual memory. The `mmap()` function does the “memory mapping”, establishing a mapping between a program's (virtual) memory address space and a file stored on your hard drive -- we call the file a “memory-mapped” file. Since memory-mapped files are viewed as a sequence of bytes (i.e., a binary file), your program needs to know how to convert bytes to and from numbers (e.g., integers). `struct` supports such conversions via [“packing” and “unpacking”](#), using format specifiers that represent the desired [endianness](#) and data type to convert to/from.

Task 1.1 Setup Pypy

Install PyPy, which is a Just-In-Time compilation runtime for python, which supports fast packing and unpacking. (As mentioned in class, C++ and Java are generally speedier than Python. However, [several projects aim to boost Python speed](#). PyPy is one of them.)

Ubuntu	<code>sudo apt-get install pypy</code>
--------	--

MacOS	Install Homebrew Run <code>brew install pypy</code>
Windows	Download the package and then install it.

Now, run the following code in the Task1 directory to learn more about the helper utility that we have provided to you for this task.

```
$ pypy task1_utils.py --help
```

Task 1.2 Warm Up (15 pts)

Get started with memory mapping concepts using the code-based tutorial in `warmup.py`.

You should study the code and modify parts of it as instructed in the file. You can run the tutorial code as-is (without any modifications) to test how it works. The warmup code is setup to pack the integers from 0 to 63 into a binary file, and unpack it back into a memory map object. You will need to modify this code to do the same thing for all odd integers in the range of 1 to 42. The lines that need to be updated are clearly marked. **You must not modify any other parts of the code.**

When you're done, you can run the following command to test whether it works as expected:

```
$ python task1_utils.py test_warmup out_warmup.bin
```

It prints `True` if the binary file created after running `warmup.py` contains the expected output.

Task 1.3 Implementing and running PageRank (25 pts)

You will implement the PageRank algorithm, using the power iteration method, and run it on the [LiveJournal dataset](#) (an online community with millions of users to maintain journals and blogs). You may want to revisit the [MMap lecture slides](#) (slide 6, 7) to refresh your memory about the PageRank algorithm and the data structures and files that you may need to memory-map. (For more details, read the [MMap](#) paper.) You will perform three steps (subtasks) as described below.

Step 1: Download the [LiveJournal graph dataset](#) (an edge list file)

The LiveJournal graph has almost 70 million edges. It is available on the [SNAP website](#). We are hosting the graph on our course homepage, to avoid high traffic bombarding their site.

Step 2: Convert the graph's edge list to binary files (you only need to do this once)

Since memory mapping works with binary files, you will convert the graph's edge list into its binary format by running the following command at the terminal/command prompt:

```
$ python task1_utils.py convert <path-to-edgelist.txt>
```

Example:

Consider the following `toy-graph.txt`, which contains 7 edges:

```
0 1
1 0
1 2
2 1
3 4
4 5
5 2
```

To convert the graph to its binary format, you will type:

```
$ python task1_utils.py convert toy-graph/toy-graph.txt
```

This generates 3 files:

`toy-graph/`

`toy-graph.bin`

binary file containing edges (source, target) in little-endian "int" C type

`toy-graph.idx`

binary file containing (node, degree) in little-endian "long long" C type

`toy-graph.json`

metadata about the conversion process (required to run pagerank)

In `toy-graph.bin` we have,

```
0000 0000 0100 0000    # 0 1    (in little-endian "int" C type)
0100 0000 0000 0000    # 1 0
0100 0000 0200 0000    # 1 2
0200 0000 0100 0000    # 2 1
0300 0000 0400 0000    # 3 4
0400 0000 0500 0000    # 4 5
0500 0000 0200 0000    # 5 2
ffff ffff ffff ffff
ffff ffff ffff ffff
ffff ffff ffff ffff
ffff ffff ffff ffff
ffff ffff ffff ffff
```

In `toy-graph.idx` we have,

```
0000 0000 0000 0000 0100 0000 0000 0000    # 0 1 (in little-endian "long long" C type )
0100 0000 0000 0000 0200 0000 0000 0000    # 1 2
...
ffff ffff ffff ffff ffff ffff ffff ffff
```

Note: there are extra values of -1 (`ffff ffff` or `ffff ffff f fff ffff`) added at the end of the binary file as padding to ensure that the code will not break in case you try to read a value greater than the file size. You can ignore these values as they will not affect your code.

Step 3: Implement and run the PageRank algorithm on the LiveJournal graph's binary files

Follow the instructions in `pagerank.py` to implement the PageRank algorithm.

You will only need to write/modify a few lines of code.

Run the following command to execute your pagerank implementation:

```
$ pypy task1_utils.py pagerank <path to JSON file for LiveJournal>
```

This will output the 10 nodes with the highest pagerank scores.

For example:

```
$ pypy task1_utils.py pagerank toy-graph/toy-graph.json
```

```
node_id score
1      0.4106875
2      0.2542078125
0      0.1995421875
5      0.0643125
4      0.04625
3      0.025
```

(Note that only 6 nodes are printed here since the toy graph only has 6 nodes.)

Copy the output for the top 10 nodes for the LiveJournal graph into **pagerank_nodes_n.txt** for n=10, 20, 30 iterations (try the `--iterations n` argument in command above). Each line in the output must contain a node and its pagerank score as tab separated values (no header needed).

You may notice that while the top nodes' ordering starts to stabilize as you run more iterations, the nodes' PageRank scores may still change. The speed at which the PageRank scores converge depends on the PageRank vector's initial values. The closer the initial values are to the actual pagerank scores, the faster the convergence.

Deliverables

- **warmup.py [8pt]**: your modified implementation as instructed in the code file
- **out_warmup.bin [5pt]**: the binary file created from your modified warmup.py
- **out_warmup_bytes.txt [2pt]**: the text file with the number of bytes calculated by you for the warm-up task. This should be automatically generated by warmup.py
- **pagerank.py [18pt]**: the python code with updated parameters
- **pagerank_nodes_n.txt [6pt]**: the top 10 node IDs and their pageranks for n iterations.
 - **pagerank_nodes_10.txt** for n=10
 - **pagerank_nodes_20.txt** for n= 20
 - **pagerank_nodes_30.txt** for n= 30
- **pagerank_time.txt [1pt]**: Report the time taken for your implementation for n=10,20,30.

Task 2: Random Forest (40 points)

You will implement a random forest classifier in Python. The performance of the classifier will be evaluated via the out-of-bag (OOB) error estimate, using a provided dataset. To refresh your memory about random forest and OOB, see Chapter 15 in the "[Elements of Statistical Learning](#)" book, [lecture](#)

[slides](#), and a nice online [discussion](#). (Here is a [blog post](#) that introduces random forests in a fun way, in layman's terms.) **You must not use existing machine learning or random forest libraries.**

You will use a wine dataset, which is often used for evaluating classification algorithms. The classification task is to determine whether the quality of a particular wine is over 7. We have mapped the wine quality scores to the binary classes of 0 and 1 for you. Wine scores from 0 to 6 (inclusive) are mapped to 0, wine scores of 7 and above are mapped to 1. You will perform binary classification on the dataset, extracted from <https://archive.ics.uci.edu/ml/datasets/Wine+Quality>. The data is stored in a comma-separated file (csv). Each line describes a wine using 12 columns: the first 11 describe the wine's characteristics, and the last column is the ground truth label for the wine quality (0/1). Note that the last column should NOT be treated as an attribute.

A. Implementing Random Forest (25 pt)

The main parameters in a random forest are:

- Which attributes of the whole set of attributes do you select to find a split?
- When do you stop splitting leaf nodes?
 - You may even choose to use decision stumps which are just 1 node deep.
- How many trees should the forest contain?

In your implementation, you may apply any variations that you like (e.g., using entropy, Gini index, or other measures; binary split or multi-way split). However, you must explain your approaches and their effects on the classification performance in a text file ***description.txt***.

We have prepared starter code written in Python (*RandomForest.py*) which you will be using. This would help you setup the environment (loading the data and evaluating your model).

You may use heuristics to improve the performance of your random forest, which you should mention in ***description.txt***. It is OK to build on top of your initial implementation and only submit the best (final) version of your random forest.

B. Computing and reporting out-of-bag error estimates (15 pt)

In random forests, it is not necessary to perform explicit cross-validation or use a separate test set for performance evaluation (also discussed in [class](#)). Out-of-bag (OOB) error estimate has shown to be reasonably accurate and unbiased. Below, we summarize the key points about OOB described in the [original article by Breiman and Cutler](#).

Each tree in the forest is constructed using a different bootstrap sample from the original data (usually, a bootstrap sample has the [same size](#) as the original dataset). Statistically, about one-third of the cases are left out of the bootstrap sample and not used in the construction of the k th tree. For each record left out in the construction of the k th tree, it can be assigned a class by the k th tree. As a result, each record will have a "test set" classification by the subset of trees that treat the record as an out-of-bag sample. The majority vote for that record will be its predicted class. The proportion of times that a predicted class is not equal to the true class of a record averaged over all records is the OOB

error estimate.

Modify the code template to compute the OOB error estimate. Report the estimate of your implementation.

Deliverables

1. **RandomForest.py:**
 - The source code of your program,
 - with detailed comments for your code.
2. **description.txt:**
 - Specific the implementation steps of the random forest and why you choose the specific approach (<75 words)
 - Report the OOB estimate

Task 3: Using Weka (20 points)

You will use [Weka](#), a popular machine learning software, to train classifiers for the same dataset used in Task 2, and to compare the performance of your random forest implementation with Weka's.

Download and install [Weka](#). Note that Weka requires Java Runtime Environment (JRE) to run. We suggest that you install the [latest JRE](#), to avoid Java or runtime-related issues.

How to use Weka:

- Load data into *Weka Explorer*: Weka supports file formats such as arff, csv, xls.
- Preprocessing: you can view your data, select attributes, and apply filters.
- Classify: under *Classifier* you can select the different classifiers that Weka offers. You can adjust the input parameters of many models by clicking on the text to the right of the *Choose* button in the Classifier section.

The above are some Weka fundamentals. There are numerous online tutorials.

A. Experiment (10 pt)

Run the following experiments. After each experiment, report your **parameters, running time, confusion matrix, and prediction accuracy**. An example is provided below, under the "Deliverables" section. For the Test options, choose **10-fold cross validation**

1. **Random Forest.** Under *classifiers* -> *trees*, select RandomForest. You might have to preprocess the data before using this classifier. (5 pt)
2. **Your choice** -- choose any classifier you like from the numerous classifiers Weka provides. You can use package manager to install the ones you need. (5 pt)

B. Discussion (10 pt)

1. Compare the Random Forest result from A1 to your implementation in Task 2 and discuss possible reasons for the difference in performance. (< 50 words, 5 pt)
2. Compare and explain the two approaches' classification results in Section A, specifically their running times, accuracies, and confusion matrices. If you have changed/tuned any of the parameters, briefly explain what you have done and why they improve the prediction accuracy. (< 100 words, 5 pt)

Deliverables

report.txt - a text file containing the Weka result and your discussion for all questions above. For example:

Section A

1.

J48 -C 0.25 -M 2

Time taken to build model: 3.73 seconds

Overall accuracy: 86.0675 %

Confusion Matrix:

	a	b	<-- classified as
33273	2079		a = no
4401	6757		b = yes

2.

...

Section B

1. The result of Weka is 86.1% compared to my result <accuracy> because...
2. I choose <classifier> which is <algorithm>...

...

Submission Guidelines

Submit the deliverables as a single zip file named **hw4-LastName-FirstName.zip** (should start with lowercase hw4). Write down the name(s) of any students you have collaborated with on this assignment, using the text box on the T-Square submission page.

The zip file's directory structure must exactly be (when unzipped):

```
hw4-LastName-FirstName/  
  Task1/  
    out_warmup.bin  
    out_warmup_bytes.txt  
    pagerank.py  
    pagerank_nodes_10.txt  
    pagerank_nodes_20.txt  
    pagerank_nodes_30.txt  
    pagerank_time.txt  
    warmup.py  
    toy-graph/... <optional>  
  
  Task2/  
    description.txt  
    hw4-data.csv  
    RandomForest.py  
  
  Task3/  
    report.txt
```

You must follow the naming convention specified above.