

Carnegie Mellon University
15-826 – Multimedia Databases and Data Mining
Fall 2013, C. Faloutsos
Homework 2, Due Date: Nov 5th, at classroom 1:30pm
Prepared by: Alex Beutel (Q4, Q5) & Vagelis
Papalexakis (Q1, Q2, Q3)

Reminders

- All homeworks are to be done **INDIVIDUALLY**.
- All written answers should be **TYPED**.
- For code submission to blackboard, please follow the template of the `hw.zip` file that you can find download from <http://www.cs.cmu.edu/~christos/courses/826.F13/HOMEWORKS/HW2/hw2.zip>. You should rename the final .zip file to [andrew-id].zip.
- FYI: Expected effort for this homework (approximate times):
 - Q1: 1-2 hours (**To be graded by Seunghak Lee**)
 - * 1-2 hours to solve and write it up.
 - Q2: 4 hours (**To be graded by Vagelis Papalexakis**)
 - * 1-2 hours to download the datasets and the software package.
 - * 0.5-1 hours to see how the software works.
 - * 1-2 hours to generate the plots and write up the solution.
 - Q3: 10 hours (**To be graded by Vagelis Papalexakis**)
 - * 3-4 hours to write the code for the 80/20 multifractal.
 - * 2-3 hours to modify the above code for the 70/20/10 multifractal.
 - * 1 hour to generate the correlation integral plots using FDNQ.
 - * 2-3 hours to code up the exact correlation integral and finalize the rest of the question.
 - Q4: 20 hours (**To be graded by Seunghak Lee**)
 - * 15 hours to code the string edit distance script
 - * 4 hours to test it for corner cases
 - * 1 hour to calculate the answers for the question
 - Q5: 10 hours (**To be graded by Alex Beutel**)
 - * 7 hours for the SQL statement
 - * 1 hour for plotting the degree distribution
 - * 2 hour for answering questions around speed

Q1 – Density Paradox [5 pts]

(On separate page)

Problem Description: In this problem we will see the density paradox. You are given $N=1$ million points on the 45-degree line (i.e., $x = y$ line), which starts from $(-M, -M)$ and up to (M, M) , where M is a large number.

Compute the density at the origin. That is, consider the $(0, 0)$ point; we are told that the number of its neighbours $NN(r)$ for radius $r = 1$ is $NN(1) = 20$. For simplicity, we consider the L_∞ norm, which means that we have a square of radius r (and thus, of side $2r$), centered at the origin, and within this square, there are 20 points.

1. [2.5 pts] Calculate the density for the following radii: [HINT: estimate the number of neighbours, for each radius, using the formula $NN(r) = C r^{D_2}$:
 - $r = 5$
 - $r = 10$
 - $r = 20$
2. [2.5 pts] What do you observe? Is the above question well posed?

What to turn in:

- **On Paper:** Please turn in the answers to the above question. (On separate page)
- **Online:** There is no need to submit something for this question.

Q2 – Fractals [10 pts]

(On separate page)

Problem Description: In this problem you will experiment with Fractal dimension, and how to use it as a feature, in order to tell whether a dataset looks 'realistic' or not. You are given 5 different datasets with coordinates in the 2-d space. You can download an archive containing all these datasets from <http://www.cs.cmu.edu/~epapalex/15826F13/data/datasets.zip> [CMU only access - please make sure you download it using an SCS machine, or login via VPN]. Each dataset is formatted as:

x_coordinate y_coordinate

1. [5 pts] For each one of those datasets, plot the correlation integral. You may use the code of the FDNQ package http://www.cs.cmu.edu/~christos/SRC/fdnq_h.zip
2. [5 pts] Using the information provided by the correlation integral, find out which one(s) (if any) of the datasets are probably not real, and justify your solution based on the correlation integral. [HINT: The fractal dimension of the fake datasets will be close to integer]

What to turn in:

- **On Paper:** (On separate page). Please turn in printouts of the correlation integrals, the slope for each of the correlation integrals, as well as the numbers of the datasets that you identified as fake.
- **Online:** Please turn in the code that you used in order to generate the required plots, including the `run.sh` bash script we have given you in the homework template, after you fill in the calls to your code. In order to execute this script, you can type `bash run.sh`. You should also turn in the plots in .pdf form.

Q3— Self Similar Time Sequence & Multifractals [25 pts]

(On separate page)

Problem Description: In this problem you will experiment with Multifractals, and the generation of bursty, realistic, self-similar time sequences, as we saw in class. Each value of the histogram you will create corresponds to, say, the number of disk accesses at a time interval.

1. [7 pts] Write code that generates and plots a multifractal bursty time-sequence, using the 80/20 rule. As we said in the course foils, a bias $b=0.8$ means that within a given time interval, 80% of the accesses happen in the right half of the interval and the remaining 20% in the left half; and this splitting of accesses happens recursively for each of those halves. You should use 2048 time-ticks, and 2048 disk accesses. [HINT: See [Wang et al.] <http://www.pdl.cs.cmu.edu/PDL-FTP/Workload/CMU-CS-01-101.pdf>]
2. [7 pts] Write code that generates and plots a multifractal bursty time-sequence, using the 70/20/10 rule. That is, instead of splitting each time interval into two parts, we split it into three; then, with probability 70% we are generating an access on the first third, 20% on the second third, and 10% on the last third. Use the same configuration for the disk accesses, as well as the same convention for choosing the segments (70% right, 20% middle, 10% left). For the number of time-ticks, you have to choose a power of 3, so pick 2187.
3. [5 pts] Plot the correlation integral for the time sequences you generated, using the FDNQ package http://www.cs.cmu.edu/~christos/SRC/fdnq_h.zip. What is the slope?
4. [5 pts] Write code that computes and plots the correlation integral using the exact, quadratic method that we saw in class. You may use your favourite language, as long as it compiles on an Andrew machine. What is the slope? Does it differ significantly from the one you found using FDNQ?

5. [1 pt] Using the correlation integral and its slope, verify that the fractal dimension D for the $b = 80$ multifractal is equal to

$$D = -\log_2(b^2 + (1 - b)^2)$$

[HINT: It is fine if the two numbers are approximately equal.]

What to turn in:

- **On Paper:** (On separate page). Please turn in (a) the numerical answers to all questions above, (b) the printouts of all the plots and (c) a printout of the code that you wrote.
- **Online:** Please turn in the code that you wrote/used in order to generate the required plots including the `run.sh` bash script we have given you in the homework template, after you fill in the calls to your code. In order to execute this script, you can type `bash run.sh`. You should also turn in the plots in .pdf form.

Q4– String Editing Distance [30 pts]

(On separate page)

Problem Description: Write a program to compute the string edit distance and the path that the string editing takes, as can be seen in the example below.
Compute the following:

1. [15 pts] Compute the string edit distance between the two strings below:

The quick brown fox jumps over the lazy dog

and

Das quik brown foxxx jumps over the lay-z dogg

Use a deletion/insertion cost of 1 and a substitution cost of 0.5.

2. [15 pts] Run the same command with deletion/insertion cost of 0.5 and substitution cost of 5.

Example output: For “some string” and “somestrng 2” with weights 1 and 0.5:

Cost: 3.5

```
s -> s
o -> o
m -> m
e -> e
  -> *
s -> s
t -> t
r -> r
i -> r
n -> n
g -> g
* ->
* -> 2
```

Note: Use an asterisk * as a placeholder when there is an insertion or deletion as in the example above.

What to turn in:

- **On Paper:** (On separate page). Please turn in the cost and the “path” that the minimum string edit distance takes as seen in the example above. Also turn *all* code you used to generate your output.
- **Online:** Please turn in your code as well as a bash script named `stredit.sh` such that the command:
`bash stredit.sh "some string" "somesrrng 2" 1 0.5`
will generate the desired output as in the example above. A template can be found at <http://www.cs.cmu.edu/~christos/courses/826.F13/HOMEWORKS/HW2/Q4/q4.tar.gz>

Note, if your program requires compiling please submit the source code and have your bash script compile your code before running it.

Q5– Graph Mining and Anomaly Detection [30 pts]

(On separate page)

Here you will learn to mine data from the graph structure. The directed graph you will be operating can be downloaded from <http://cs.cmu.edu/~abeutel/patents.csv>. In this file each line is an edge specified as `fromNode,toNode`. Please complete the following tasks using SQL.

Problem Description:

1. [3 pts] How many nodes are there in the graph? (Denote this as N .)
2. [3 pts] What are the total number of edges (or 1 hop away neighbors) in the graph? How long does it take to find this value?
3. [5 pts] Find and plot the out degree distribution for the graph. Give the plot.
4. [3 pts] At what degree does there appear to be a surprising number of nodes with that out degree?

For the following questions, consider the graph to be undirected. The undirected dataset can be downloaded from http://cs.cmu.edu/~abeutel/undirected_patents.csv:

5. [3 pts] What is the average out degree of each node? (Denote this as d_{avg} .)
6. [3 pts] What is the maximum out degree? (Denote this as d_{max} .)
7. [3 pts] **(Note: If your query takes over an hour for this question, stop it.)** What is the total number of 2 step away neighbors for all nodes?

Also, answer the following related questions:

- (a) [4 pts] Select which of the below statements about the total number of 2 step away neighbors are true.
 - i. The value is approximately $N \cdot d_{\text{avg}} \cdot d_{\text{avg}}$ because everybody has approximately d_{avg} neighbors who have approximately d_{avg} neighbors.
 - ii. The value is more than $N \cdot d_{\text{avg}}^2$
 - iii. The value is at most $\sum_{i=1}^N d_i^2$ where d_i is the degree of node i .
 - iv. The value is at least $\sum_{i=1}^N d_i^2$ where d_i is the degree of node i .
 - v. The value is more than d_{max}^2
- (b) [3 pts] Calculate the values of the following using SQL:
 - i. $N \cdot d_{\text{avg}}^2$
 - ii. $\sum_{i=1}^N d_i^2$
 - iii. d_{max}^2

What to turn in:

- **On Paper:** (On separate page). Please turn in the answers to all questions above as well as all SQL commands needed for questions #1,2,3,5,6,7.
- **Online:** Please turn in your SQL code and a bash script to generate the answers to all of the questions above. You can follow the template bash script from <http://www.cs.cmu.edu/~christos/courses/826.F13/HOMEWORKS/HW2/Q5/q5.tar.gz>