RENSSELAER POLYTECHNIC INSTITUTE

MACHINE LEARNING AND OPTIMIZATION
CSCI 4961/6961     SECTION: 01

# Triangle Counting Problem Set
## Answer Key

*Antonia Calia-Bogan*
*Richie Massimilla*
*Gabriel Orlanski*

taught by
Prof. Alex GITTENS

December 2, 2020

**Question 1.** (*Asymptotic bound on spectral counting*) Let $A$ be an adjacency matrix for an undirected graph $G$. Show that counting triangles in $G$ using $\triangle(G) = \frac{1}{6}\mathrm{tr}(A^3)$ has asymptotically bound by $O(|E||V|)$ where $E$ and $V$ are the sets of edges and nodes in $G$ respectively.

It's recommended that you first express the cost of this process using $\deg(V_i)$, the degree of the $i^{\text{th}}$ node. The degree of a node is defined as the number of edges incident to that node. Then, you can make use of the equation $\mathbb{E}[\deg(V_i)] = |E|/|V|$ where $i$ is sampled uniformly from the integers between 1 and $|V|$ inclusive.

**Question 2.** (*Random trace estimation*) Let $A$ be an $n \times n$ symmetric matrix. Let $j$ be a uniform random integer between 1 and $n$, and let $\mathbf{z} = \mathbf{e}_j$ be the $j^{\text{th}}$ standard basis vector. Show that $n\mathbf{z}^\top A\mathbf{z}$ is an unbiased estimator of $\mathrm{tr}(A)$ by showing that the following are true

$$\mathbb{E}[n\mathbf{z}^\top A\mathbf{z}] = \mathrm{tr}(A)$$

$$\mathrm{Var}(n\mathbf{z}^\top A\mathbf{z}) = n\mathrm{tr}(A^2) - \mathrm{tr}^2(A)$$

Then argue why the TraceTriangle algorithms have lower variance than this.

**Question 3.** (*Implementing TriangleTrace$_N$*) You will implement the TriangleTrace$_N$ algorithm. The algorithm is as follows

> **Algorithm:** `TraceTriangle`$_N$
>
> **Input:** $\gamma \longleftarrow$ a scalar
> **Output:** $\triangle = $ `TraceTriangle`$_N$($G$, undirected graph with $n$ nodes)
> Form the adjacency matrix $A \in \mathbb{R}^{n \times n}$
> $M = \lceil \gamma \ln^2 n \rceil$
> **for** $i \in 1, \ldots, M$ **do**
> $\quad$ Form the vector $\mathbf{x} = [x_0, \ldots, x_n]$,
> $\qquad$ where $x_k \sim \mathcal{N}(0, 1)$ are i.i.d
> $\qquad$ and $k \in 1, \ldots, n$
> $\quad$ $y \longleftarrow A\mathbf{x}$
> $\quad$ $T_i \longleftarrow (y^T Ay)/6$
> **end**
> $\triangle \longleftarrow \frac{1}{M} \sum_{i=1}^{M} T_i$

This will require writing a function `tracetriangle(A,gamma)` which takes as input `A`, an adjacency matrix, and `gamma`, a hyperparameter that scales the number of iterations. You must implement this using a sparse matrix representation such as `scipy.sparse.csr_matrix` objects. Use any undirected graphs from the Stanford SNAP dataset (https://snap.stanford.edu/data/).

Experiment with different values of gamma and produce a graph plotting `gamma` values against mean absolute error which is computed by

$$\delta = \frac{1}{N} \sum_{i=1}^{N} \left| \frac{v_A - v_E}{v_E} \right|$$

where $N$ is the number of datasets that a particular value of `gamma` was evaluated over, $v_A$ is the ground truth triangle count, and $v_E$ is the triangle count given by `TraceTriangle`$_N$. On the following page a function to load graphs as sparse adjacency matrices and a function for computing the spectral count, which will be used as the ground truth counts, are given.

```python
1   # IMPORT NX
2   import networkx as nx
3
4   # GRPAH LOADING FUNCTION
5   def load_graph(path, data=True, delim=None):
6       """
7       Given the path to a csv file containing a row for every edge,
8       parse the data into an adjacency matrix. Each row should have two
9       elements, one for each node in the edge.
10
11      Parameters
12      ------------------------------------------------------------------
13      path : string
14      A path to the csv file for the graph.
15      data : list of pairs
16      Tuples specifying dictionary key names and types for edge
17      data.
18      delim : string
19      Delimiter string for graph read.
20
21      Returns
22      ------------------------------------------------------------------
23      out : csr_matrix
24      The graph's adjacency matrix.
25      """
26      with open(path, 'rb') as f:
27          G = nx.read_edgelist(f, data=data, delimiter=delim)
28      # The adjacency list is returned as a csr_matrix as a computational
29      # time improvement since most real graphs will be extremely sparse.
30      # This turns |V|^2 operations into |E| operations which is a huge
31      # improvement.
32      A = nx.to_scipy_sparse_matrix(G)
33      return A
34
35  # SPECTRAL COUNT
36  def gt_count(A):
37      """
38      Uses spectral counting to calculate the exact total number of
39      triangles in a graph from its adjaceny matrix.
40
41      Parameters
42      ------------------------------------------------------------------
43      A : csr_matrix
44      Adjacency matrix of the graph.
45
46      Returns
47      ------------------------------------------------------------------
48      out : int
49      Exact total count of triangles in the graph.
50      """
51      cubed = A ** 3
52      trace = cubed.diagonal().sum()
53      return trace // 6  # This will be an integer regardless
```