

Electronics and Computer Science
Faculty of Physical Sciences and Engineering
University of Southampton

Vlad Sebastian Velici
December 8, 2014

Similar nodes in large graphs

Project Supervisor: Dr. Adam Prügel-Bennett
Second Examiner: Dr. Sasan Mahmoodi

A progress report submitted for the award of
MEng Computer Science with Artificial Intelligence

Contents

Abstract	2
Introduction	3
Brute-Force Approach	3
Approximative Approach	4
Vectorisation	7
Implementation	8
Evaluation	9
Comparison of approximative and brute-force approaches	9
Visualisation	10
Discussion of further evaluation	11
Limitations of the Current Implementation	12
References	13

ABSTRACT

Given a large graph, in the range of millions or billions of nodes, it is difficult to efficiently find similar nodes. Examples of applications of finding similar nodes in data represented as graphs are all around the web nowadays: *people you may know* and *who to follow* suggestions on social networks, or topics that may interest you on different websites.

This project will explore ways to obtain meaningful information of this nature from large datasets represented as graphs.

The goals of this project are developing efficient ways to find similar nodes in large graphs for different use cases and implementing an open-source software library to provide everyone with this functionality.

INTRODUCTION

Plenty of datasets are or can be represented as graphs where vertices represent entities and edges represent relationships between entities. A problem of interest is to find entities that are similarly connected. Example instances of this problem are finding *people you may know* in a social network, people with common interests from research publications repositories or identifying possible duplicates in a dataset.

It is not computationally feasible to calculate the exact similarities between vertices of very large graphs. In this project, we investigate a method to compute an approximation of similarities between vertices and attempt to evaluate its performance (accuracy) on different datasets.

BRUTE-FORCE APPROACH

The similarity algorithm is described regarding the undirected graph $G(v, \epsilon)$ with adjacency matrix M . The number of nodes is $N = |v|$ and the number of edges is $|\epsilon|$. In order to derive an algorithm to find similarities between vertices in a graph, a measure of similarity between two vertices must be defined. M is the adjacency matrix of the graph and $M_{ij} = [(i, j) \in \epsilon]$, that is $M_{ij} = 1$ if there is an edge between vertices i and j , otherwise $M_{ij} = 0$. Let D be a normalised adjacency matrix such that each column sums up to 1:

$$D_{ij} = \frac{1}{\sum_{j=1}^N M_{ij}} \quad (1)$$

Let δ_i be the i^{th} column of the $N \times N$ identity matrix, and \mathbf{C}_i a vector that represents how vertex i is connected to the graph. In a brute-force approach, \mathbf{C}_i can be computed as:

$$\mathbf{C}_i = \sum_{t=0}^{\infty} \mu^t D^t \delta_i \quad (2)$$

Where $\mu \in (0, 1)$ is a penalising factor. μ controls how much the far connections are considered into the similarities. Therefore, a large μ enforces consideration of very far connections whereas a small μ considers far connections less. The similarity between two vertices i and j

is the (squared) norm of the difference between \mathbf{C}_i and \mathbf{C}_j :

$$\|\mathbf{C}_i - \mathbf{C}_j\|^2 \quad (3)$$

The brute-force approach was run on small randomly generated connected graphs (200 vertices) with 50 iterations and it gave meaningful results but it does not scale for large datasets. In the next section an algorithm that scales up to large datasets is described.

APPROXIMATIVE APPROACH

Geometrically, if vectors \mathbf{C}_i and \mathbf{C}_j start from the origin, then \mathbf{C}_i , \mathbf{C}_j and $\mathbf{C}_i - \mathbf{C}_j$ describe a triangle. Let θ be the angle between \mathbf{C}_i and \mathbf{C}_j . From the Law of Cosines:

$$\|\mathbf{C}_i - \mathbf{C}_j\|^2 = \|\mathbf{C}_i\|^2 + \|\mathbf{C}_j\|^2 - 2\|\mathbf{C}_i\|\|\mathbf{C}_j\|\cos\theta \quad (4)$$

Take the dot product $\mathbf{C}_i^T \mathbf{C}_j = \|\mathbf{C}_i\|\|\mathbf{C}_j\|\cos\theta$. The angle between \mathbf{C}_i and itself is 0 and $\cos 0 = 1$, thus $\|\mathbf{C}_i\|^2 = \mathbf{C}_i^T \mathbf{C}_i$. Now *Equation 4* becomes:

$$\|\mathbf{C}_i - \mathbf{C}_j\|^2 = \mathbf{C}_i^T \mathbf{C}_i + \mathbf{C}_j^T \mathbf{C}_j - 2\mathbf{C}_i^T \mathbf{C}_j \quad (5)$$

Therefore, it is sufficient to be able to compute all dot products $\mathbf{C}_i^T \mathbf{C}_j$ (for all $i, j \in v$) to compute the similarities between all vertices of the graph. Computing and storing all $\mathbf{C}_i^T \mathbf{C}_j$; $i, j \in v$ takes too much time and memory for a very large dataset if we use the brute-force approach.

An approximative algorithm to efficiently compute the dot products $\mathbf{C}_i^T \mathbf{C}_j$ will be described. Let W be a diagonal matrix:

$$W_{ij} = \frac{[i=j]}{\sum_{j=1}^N M_{ij}} \quad (6)$$

Now *Equation 1* can be written as $D = WM$, thus

$$D^t = (WM)^t \quad (7)$$

Define a new matrix A :

$$A = W^{1/2} M W^{1/2} \quad (8)$$

From *Equation 7* and *Equation 8*:

$$D^t = W^{1/2} A^t W^{-1/2} \quad (9)$$

The matrices $W^{1/2}$ and $W^{-1/2}$ are easy to compute because W is a diagonal matrix. Apply the matrix operation only on the diagonal elements and obtain $(W^{1/2})_{ii} = (W_{ii})^{1/2}$ and $(W^{-1/2})_{ii} = (W_{ii})^{-1/2}$.

Let λ_i be the i^{th} eigenvalue and $V^{(i)}$ be the i^{th} eigenvector ($\forall i \in \{1, 2, \dots, N\}$). V is a matrix of all computed eigenvectors (such that $V^{(i)}$ is the i^{th} column of V), and Λ is a diagonal matrix of all computed eigenvalues:

$$V = \begin{bmatrix} | & | & & | \\ V^{(1)} & V^{(2)} & \dots & V^{(N)} \\ | & | & & | \end{bmatrix} \quad (10)$$

$$\Lambda = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & & \lambda_N \end{bmatrix} \quad (11)$$

The eigendecomposition of the matrix A is $A = V \Lambda V^{-1}$. The graph G is an undirected graph which makes A a real symmetric matrix, thus $V^{-1} = V^T$. The eigendecomposition of the matrix A is therefore:

$$A = V \Lambda V^T \quad (12)$$

$$\Rightarrow A^t = (V \Lambda V^T)^t = V \Lambda^t V^T \quad (13)$$

$$\Leftrightarrow A^t = \sum_{a=1}^N \lambda_a^t V^{(a)} V^{(a)T} \quad (14)$$

From *Equation 9* and *Equation 14*:

$$D^t = W^{1/2} \sum_{a=1}^N \lambda_a^t V^{(a)} V^{(a)T} W^{-1/2} \quad (15)$$

Plug D^t from *Equation 15* into *Equation 2*:

$$C_i = \sum_{t=0}^{\infty} \mu^t W^{1/2} \sum_{a=1}^N \lambda_a^t V^{(a)} V^{(a)T} W^{-1/2} \delta_i \quad (16)$$

Rearrange the equation and obtain:

$$\mathbf{C}_i = W^{1/2} \sum_{a=1}^N \sum_{t=0}^{\infty} \mu^t \lambda_a^t \mathbf{V}^{(a)} \mathbf{V}^{(a)T} W^{-1/2} \delta_i \quad (17)$$

Note that $\sum_{t=0}^{\infty} \mu^t \lambda_a^t$ is a sum of a geometric series with the first term 1 and ratio $\mu \lambda_a$, therefore:

$$\sum_{t=0}^{\infty} \mu^t \lambda_a^t = \frac{1}{1 - \mu \lambda_a} \quad (18)$$

Observe $\mathbf{V}^{(a)T} W^{-1/2} \delta_i$ is a scalar:

$$\mathbf{V}^{(a)T} W^{-1/2} \delta_i = \mathbf{V}_i^{(a)} W_{ii}^{-1/2} \quad (19)$$

Substitute *Equation 19* and *Equation 18* back into *Equation 17* and obtain:

$$\mathbf{C}_i = W^{1/2} \sum_{a=1}^N \frac{1}{1 - \mu \lambda_a} \mathbf{V}^{(a)} \mathbf{V}_i^{(a)} W_{ii}^{-1/2} \quad (20)$$

An approximation of the vectors $\mathbf{C}_i (\forall i \in v)$ can be obtained by only using the m leading eigenvectors and eigenvalues of A instead of all of them. Compute the largest m eigenvalues and eigenvectors of the matrix A using the Lanczos Method for real symmetric matrices[1]. For directed graphs the Arnoldi Iteration[2] will be used instead. These algorithms perform well on large sparse matrices. In real datasets it is very often the case that the number of edges is orders of magnitude smaller than the number of vertices squared, resulting in sparse adjacency matrices and a sparse A . Define $\hat{\mathbf{C}}_i$ to be an approximation of \mathbf{C}_i :

$$\hat{\mathbf{C}}_i = W^{1/2} \sum_{a=1}^m \frac{1}{1 - \mu \lambda_a} \mathbf{V}^{(a)} \mathbf{V}_i^{(a)} W_{ii}^{-1/2} \quad (21)$$

To simplify the equation, let Z be a $N \times m$ matrix such that:

$$Z_{ij} = \frac{\mathbf{V}_i^{(j)} W_{ii}^{-1/2}}{1 - \mu \lambda_j} \quad (22)$$

Substitute in *Equation 21* and get:

$$\hat{\mathbf{C}}_i = W^{1/2} \sum_{a=1}^m Z_{ia} \mathbf{V}^{(a)} \quad (23)$$

The dot product between $\hat{\mathbf{C}}_i$ and $\hat{\mathbf{C}}_j$ becomes:

$$\hat{\mathbf{C}}_i^T \hat{\mathbf{C}}_j = \left(W^{1/2} \sum_{a=1}^m Z_{ia} \mathbf{V}^{(a)} \right)^T \left(W^{1/2} \sum_{a'=1}^m Z_{ja'} \mathbf{V}^{(a')} \right) \quad (24)$$

$$\Leftrightarrow \hat{\mathbf{C}}_i^T \hat{\mathbf{C}}_j = \sum_{a=1}^m Z_{ia} \mathbf{V}^{(a)T} W \sum_{a'=1}^m Z_{ja'} \mathbf{V}^{(a')} \quad (25)$$

$$\Leftrightarrow \hat{\mathbf{C}}_i^T \hat{\mathbf{C}}_j = \sum_{a=1}^m \sum_{a'=1}^m Z_{ia} Z_{ja'} \mathbf{V}^{(a)T} W \mathbf{V}^{(a')} \quad (26)$$

Define a new $m \times m$ matrix Q such that:

$$Q_{ij} = \mathbf{V}^{(i)T} W \mathbf{V}^{(j)} \quad (27)$$

\mathbf{Z}_i is the i^{th} row of the matrix Z . Substitute Q into Equation 26 and obtain:

$$\hat{\mathbf{C}}_i^T \hat{\mathbf{C}}_j = \sum_{a=1}^m \sum_{a'=1}^m Z_{ia} Z_{ja'} Q_{aa'} \quad (28)$$

$$\Leftrightarrow \hat{\mathbf{C}}_i^T \hat{\mathbf{C}}_j = \mathbf{Z}_i^T Q \mathbf{Z}_j \quad (29)$$

From Equation 29 it is deduced that it is only required to store the matrixes Q and Z , which only takes $O(m^2 + Nm)$ memory, instead of $O(N^2)$ to store all $\hat{\mathbf{C}}_i$ vectors.

To compute the similarity between two vertices, i and j , use the norm formula from Equation 29 and obtain:

$$\|\hat{\mathbf{C}}_i - \hat{\mathbf{C}}_j\|^2 = \mathbf{Z}_i^T Q \mathbf{Z}_i + \mathbf{Z}_j^T Q \mathbf{Z}_j - 2 \mathbf{Z}_i^T Q \mathbf{Z}_j \quad (30)$$

After computing the matrices Z and Q , the time complexity of computing the similarity between two nodes is the same as multiplying three small matrices (of sizes $1 \times m$, $m \times m$ and $m \times 1$).

Vectorisation

It is usually more efficient to have vectorised implementations (to use matrix operations in favour of iterative solutions) in programming languages like Matlab, Octave or R. In this subsection is presented a way of computing the matrices Z and Q using matrix operations. For this purpose, in this subsection, \mathbf{V} only contains the computed eigenvectors (is of size $N \times m$), and $\mathbf{\Lambda}$ only contains the computed eigenvalues (is of size $m \times m$).

To vectorise the computation of Z from *Equation 22*, define a new diagonal matrix R :

$$R = \begin{bmatrix} \frac{1}{1-\mu\lambda_1} & 0 & \dots & 0 \\ 0 & \frac{1}{1-\mu\lambda_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \\ 0 & 0 & \dots & \frac{1}{1-\mu\lambda_m} \end{bmatrix} \quad (31)$$

And then define Z as:

$$Z = W^{-1/2} \bigvee R \quad (32)$$

A vectorised version of *Equation 27* (to compute Q):

$$Q = \bigvee^T W \bigvee \quad (33)$$

IMPLEMENTATION

A choice had to be made on what tools and languages to use for the implementation of this algorithm. It involves linear algebra so programming languages that support such operations or have good libraries for this purpose were candidates. Matlab was chosen for this task because it has very strong linear algebra features built-in, it is easy and fast to try different experiments (can also be used interactively) and programs are usually smaller in size (lines of code) than their variants in C++ or other languages. At a later stage when the algorithm is closer to its final version, porting it to a language like C++ will be considered.

Matlab has the function $eigs(A, m)$ to compute the leading m eigenvectors and eigenvalues. Internally, it uses the Lanczos Method for Hermitian matrices and the Arnoldi Iteration for non-Hermitian matrices.[3]

For the implementation of both the approximative approach and the brute-force algorithm, sparse matrices were used. A sparse matrix in Matlab only stores the non-zero values of a matrix, therefore a sparse matrix only uses memory for the positions where it holds values.[4] For instance, to store the adjacency matrix of a very large graph $G(v, \epsilon)$, using a dense matrix would take $O(|v|^2)$ memory whereas using a sparse matrix only takes $O(|\epsilon|)$ memory.

The approximative approach was implemented in Matlab as a couple of Matlab functions.

The function signatures along with descriptions is presented in *Table 1*.

similarity(adj, mu, m)	Computes and returns the matrices Q and Z.
sim2(q,z,i,j)	Computes the similarity between vertices i and j.
simlist(q,z,i)	Computes the similarity between vertex i and all other vertices.
top(q,z,i)	Uses simlist(q,z,i) and sorts the results.
randomGraph(nodes, vertices, directed)	Generates a random connected graph.
readEprints()	Reads ePrints data from a pre-processed file.

Table 1: Implementation interface. Parameters: q represents the matrix Q, z represents the matrix Z, adj represents the adjacency matrix, and i and j represent vertex numbers.

The brute-force algorithm was implemented as a separate set of Matlab functions with a very similar interface.

EVALUATION

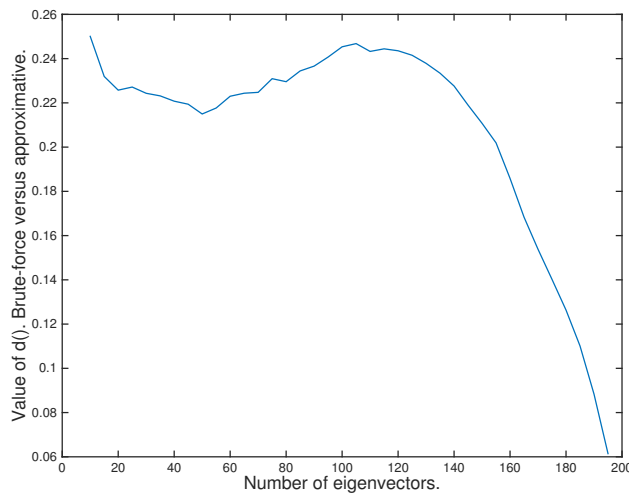
Comparison of approximative and brute-force approaches

A simple evaluation metric of how well the approximative approach retains the most important information from the graph is to evaluate how different are the top similar vertices for the brute-force approach and the approximative approach with different number of eigenvalues and eigenvectors used.

For a dataset of N vertices, let B_{ij}^m be the position of vertex j in the top similar vertices of vertex i for the brute-force approach with m iterations. Similarly let T_{ij}^m be the position of vertex j in the top similar vertices of vertex i for the approximative approach with m eigenvectors. Define the difference between two such tops to be:

$$d(B^m, T^{m'}) = \frac{1}{N^3} \sum_{i=1}^N \sum_{j=1}^N \text{abs}(B_{ij}^m - T_{ij}^{m'}) \quad (34)$$

The equation is normalised by N^3 because the maximum value we can get lies in the interval (N^2, N^3) . The results can be observed in *Figure 1*. The choice of the penalising factor $\mu = 0.5$ was used for all the runs. The difference d was evaluated between the brute-force approach with $m = 20$ iterations and the approximative approach with $m' = 5$ to 195 eigenvectors. We can observe that if we use the majority of eigenvalues, the value of $d()$ is very small. It is



(a) A plot of the differences on a range of eigenvalues.

Parameters	Value of $d()$
$d(B^{20}, T^{10})$	0.2501
$d(B^{20}, T^{20})$	0.2257
$d(B^{20}, T^{30})$	0.2243
$d(B^{20}, T^{40})$	0.2208
$d(B^{20}, T^{50})$	0.2150
$d(B^{20}, T^{100})$	0.2453
$d(B^{20}, T^{199})$	0.0325

(b) Table of some of the values.

Figure 1: Difference of results between brute-force and approximative approaches.

computationally infeasible to do that on a large dataset thus we will aim for a number of eigenvectors m' in the range of 20 to 40, where the difference seems to be at a local minimum (on this small dataset). The local minimum might differ from dataset to dataset depending mainly on the size of the dataset. For completeness, the experiment was run on a randomly generated connected graph with 500 vertices and 700 edges. The first local minimum of $d()$ was at around 100 eigenvectors with a value of approx. 0.2. At 40 and 20 eigenvectors, the values of $d()$ obtained are 0.2297 and 0.2557, respectively.

Although the evaluation metric described above might be useful to compare the two approaches directly, it does not tell anything about how well the algorithm achieves its goal of finding similar entities in a real dataset. Even if the computation of the brute-force approach would be feasible on large datasets for comparison, it would still not verify the underlying intuition of this algorithm.

Visualisation

A subjective method of evaluation on real data is attempting to visualise the algorithm output on real data and deduce whether it gives sensitive results. This method has the disadvantage of being subjective but is nonetheless a start of the journey of building an evaluation method

for the algorithm.

Discussion of further evaluation

None of the evaluation methods described above give an objective, efficient and automatic evaluation metric. In this subsection we discuss a few possible evaluation methods of the algorithm.

To objectively evaluate how well the outputs of the algorithm reflect the reality of the dataset it is being used on, we need to carefully define the goals of the algorithm in respect of the dataset. We only look for datasets where the *similarities* between entities can be obtained exclusively from the graph structure. We need labelled (known *similarities* between entities) datasets that have only one relationship type from which we can exclusively compute *similarities* between nodes.

If the dataset used for evaluation has more than one type of relationship (more features), the evaluation is subjective to the choice of the weights for each relationship type. The reference similarities (from dataset or obtained by other means) might be based on different weights of each feature of the dataset, therefore using equal weights in the graph representation will likely lead to a poorer performance and biased evaluation.

Given some datasets that fit the requirements from above, we can run our algorithm and compare the results with the reference similarities. In reality datasets have many features and a complex graph structure with relationships having different weights, therefore this evaluation method cannot be done in practice.

A similar way is to compute similarities using another algorithm (call it reference algorithm) and only one feature of a dataset. Use the same feature to build a graph structure, then compare the results of our algorithm with the results of the reference algorithm used.

Given a labelled dataset that fits the requirements from above, we can run our algorithm and compare the results with the reference similarities. A single dataset can result in biased evaluations, therefore more such datasets are required. The downside is that datasets with known similarities between entities are not very common, and because the term *similarity* is ambiguous (e.g. why is an entity *similar* to another? In respect to which features?). Given an

unlabelled dataset, we can obtain a labelled dataset using other algorithms and techniques, but the evaluation will be in respect to the labelling method used.

An online evaluation is a way to evaluate the system by getting feedback from users. In this case, we would deploy the algorithm on a experimentation website and display the results to live users, while making sure to track the way users interact with the outputs of the algorithm.

LIMITATIONS OF CURRENT IMPLEMENTATION

The current implementation of this algorithm has various limitations which are discussed along with possible improvements.

Only undirected graphs Directed graphs are not currently supported. In practice, datasets have meaningful unidirectional relationships (e.g. in a social network person A follows person B, but B does not follow A), and often datasets are represented as directed graphs rather than undirected graphs. The algorithm can be adapted to support both directed and undirected graphs but it will have the disadvantage of requiring to compute V^{-1} (for undirected graphs, $V^{-1} = V^T$).

Ignored relationship types In the real world datasets might have different types of relationships between entities. Some relationship types might be more relevant than others in finding a specific result. For instance, in a social network two people being friends might be more relevant for recommending new friends than two people following the same topic. The current implementation of the algorithm can compute similarities if it is given a weighted (symmetric) adjacency matrix but it does not have a way to automatically obtain these weights from the dataset.

Not distributed What if the dataset is too large to fit into main memory? The algorithm is currently only designed to run on one machine. Methods of distributing the algorithm on more than one machine over a network will be investigated in the future.

REFERENCES

- [1]: Calvetti, Daniela, L. Reichel, and Danny Chris Sorensen. "An implicitly restarted Lanczos method for large symmetric eigenvalue problems." *Electronic Transactions on Numerical Analysis* 2.1 (1994): 21.
- [2]: Lehoucq, Richard B., and Danny C. Sorensen. "Deflation techniques for an implicitly restarted Arnoldi iteration." *SIAM Journal on Matrix Analysis and Applications* 17.4 (1996): 789-821.
- [3]: "Matlab 2014b Documentation." Largest Eigenvalues and Eigenvectors of Matrix. The MathWorks, Inc., 3 Oct. 2014. Web. Accessed 07 Dec. 2014.
<<http://uk.mathworks.com/help/matlab/ref/eigs.html>>.
- [4]: Gilbert, John R., Cleve Moler, and Robert Schreiber. "Sparse matrices in MATLAB: design and implementation." *SIAM Journal on Matrix Analysis and Applications* 13.1 (1992): 333-356.