

MLSS-17 Tutorial: Social Network Analysis

Manuel Gomez Rodriguez, Utkarsh Upadhyay, and Isabel Valera

1 Building the package

You will need to work from a MAC or a linux machine (either locally, or remotely using ssh). Download the code package (<https://github.com/Networks-Learning/mlss-2017>) and compile it as:

```
$ cd code/  
$ make
```

Several warnings might appear when you compile the code package. As long as you do not get an Error message, you are fine. Additionally, please **download and install**:

- Graphviz: <http://www.graphviz.org/>
- Gnuplot: <http://www.gnuplot.info/>
- Gephi¹ (on your personal computer): <http://www.gephi.org/users/download/>

Now you can go to the subfolder `mlss-17/` where you will start creating networks and coding!

2 Introduction to SNAP: Network Modeling

Creating graphs and networks

We will start by creating manually two small examples of a graph and a network (*i.e.*, a graph with some data stored in the nodes and/or edges) using SNAP:

```
$ ./generate_network
```

¹If Gephi crashes in Linux, try to include the following line at the beginning of the linux gephi executable at `bin/gephi` (in the gephi package), which is actually a script: `export LIBGL_ALWAYS_SOFTWARE=1`. Alternatively, it also works by starting gephi as: `LIBGL_ALWAYS_SOFTWARE=1 ./gephi`. For more details visit <https://github.com/gephi/gephi/issues/1192>

The above program should print out the nodes and connectivity, and generate four files `graph.png`, `graph.dot`, `network.png` and `network.dot`.

Question 1.1. Do you recognize the graph and the network structures from the text print out? Draw the graph and the networks by hand. Check whether you are right by opening `graph.png` and `network.png`.

Question 1.2. Open the source code, `generate_network.cpp`, and find out how the code works. To better understand the code, you can check out a quick intro to SNAP at <http://snap.stanford.edu/snap/quick.html>.

Coding 1.1. Change `generate_network.cpp` to create a different graph and network. In particular, the code you need to modify is delimited by `TODO` comments. Once you have changed the file, compile the code by running `make` at the current folder (`mlss-17/`). After re-running `generate_network`, check whether the outputted files `graph.png` and `network.png` look like you thought.

If you create graphs and networks with a large number of nodes or edges, you will notice that `graph.png` and `network.png`, created using Graphviz, soon become messy. Gephi will help us to visualize such large graphs later.

Creating large random graphs and networks

Now, we will generate large random networks that mimic the structure of real-world social networks which we will use in the next section to simulate propagation processes or cascades. We initially consider a model of directed real-world social networks: the **Forest Fire** (scale free) model [1]. We will generate networks that have a weight associated to each edge, which we will use for our propagation models. To make it easier for you, we have coded a program (`generate_random_network`) which generates both Forest Fire and also Kronecker networks (another type of random networks).

Question 1.3 Generate several Forest Fire networks with different network parameters and compare their properties (in- and out-degree distribution, clustering coefficient, number of triangles, etc...) using `compute_properties`. In particular, you should try different burning probabilities (option `-g`), number of nodes and edge weights range (options `-la` and `-ua`; generated uniformly at random). To help you get started, check the following command line, which generates a Forest Fire network with 1,000 nodes with a burning forward probability of 0.2 and a burning backward probability of 0.17:

```
$ ./generate_random_network -t:1 -g:"0.2;0.17" -n:1000 -o:"ff-  
network" -st:1 -sg:1
```

The program will output two files: `ff-network.txt` and `ff-network.gexf`, which include a text and Gephi version of the network. Once you have generated several networks, can you tell what are the burning probabilities controlling for?

You can learn how the Forest Fire model actually works by reading Sections 4.2.1. and 4.2.2 of Leskovec et al. ([4], <http://www.cs.cmu.edu/~jure/pubs/powergrowth-tkdd.pdf>).

Coding 1.2 (optional): implement another model of directed random networks using SNAP. You will only need to edit `mlss.cpp` at the TODO comment (there is only one within the file). It may be helpful to use some of the random graph generators provided by SNAP, check out `code/snap-core/ggen.h`. You will need to call the functions writing something like `TSnap::method(Graph, ...)`. Do not forget to compile the code again by typing `make` after you finish editing.

Take into account that some of the graph generators produce an undirected graph (`PUNGraph`), and you will have to convert it to a directed graph (`PNGraph`). In those cases, decide the direction of each original undirected edge uniformly at random using `TFlt::Rnd.GetUniDev()`.

3 Introduction to Gephi: Network visualization

Layout, properties and manipulation

Use Gephi to open any of the `.gexf` network files which are located in the folder `code/mlss-17/data/`. They are networks based on real data. Additionally, you can also open the random networks you generated previously or generate new ones by running again `./generate_random_network` (use `-sg:1` to generate a Gephi file. You can also convert the networks text files (in case you did not use `-sg:1`) using `convert_to_gephi`.

Question 2.1. Visualize the networks using different layout algorithms. For this, you should use the Layout box on the left bottom corner. Which of them do you think works *better*? Can you draw any conclusions based on a particular layout? Can you guess what do the real networks represent? If you would like to know more about Gephi layouts, have a look to the tutorial <http://gephi.org/tutorials/gephi-tutorial-layouts.pdf>.

Question 2.2. Manipulate the networks. For this, use the Filters box on the right top corner. Have a look at the topology menu and try several filters (e.g., degree range, giant component, ego networks). You will need to drag and drop it the filter to the query box on the right bottom corner. Find out how to stack filters (create subfilters).

Question 2.3. Compute properties of the network. For this, go to the Window menu and select Statistics to open up the Statistics menu. Compute several statistics, including modularity. Then, partition the network and rank nodes using the Partition and Rank boxes on the top left corner; use different measurements for the partition and rank. Do you find clear partitions or clusters in the networks? Which measures work *better*? Do the networks have a hierarchical or core-periphery structure?

You can learn more about Gephi by reading the wiki at <http://wiki.gephi.org/>. You can also find an example of a research project which uses Gephi extensively at <http://snap.stanford.edu/infopath/>, where all information networks images (<http://snap.stanford.edu/infopath/graphs.html>) and videos of dynamic information networks (<http://snap.stanford.edu/infopath/videos.html>) have been generated with Gephi.

4 Information Propagation and Influence Maximization

Until here, you have learned a bit about SNAP and Gephi, and get to know two mathematical models of social and information networks. Now, we will use this knowledge to study information propagation over networks, and study the influence (spread) maximization problem. In particular:

1. You will simulate the propagation of contagions over networks using a widely known probabilistic models of influence and information propagation: the discrete time independent cascade model [2]. You can think of a *contagion* as an information unit which appears at some node of a network and then spreads like an epidemic from node to node over the edges of the underlying network. In case of information diffusion, the contagion represents a piece of information [3] and infection events correspond to times when nodes mention or copy the information from one of their neighbors in the network.
2. You will solve the influence maximization problem: you will look for the most influential source node set of a given size in a network. A contagion that starts spreading in such an influential set of nodes is expected to reach the greatest number of nodes in the network.
3. If time remains, given a fixed source set we will look for the optimal set of edges to add to a network in order to reach the greatest number of nodes in the network.

Discrete-time independent cascade model

The remainder of the practical use the discrete-time independent cascade model, and it will help you to understand it well. The model works as follows: A spreading process (or cascade) starts when a source node set A becomes infected at epoch $t = 0$. Then, source nodes have a single chance to try to infect their children (*i.e.*, neighbors that they can reach directly through an outgoing edge) at epoch $t = 1$ with some probability. When a child becomes active, it has a single chance of activating each currently inactive children at epoch $t = 2$, and so on. The activation attempts succeeds with probability p_{vw} , where v denotes parent and w denotes children. Here, note that the time is modeled only implicitly through the epochs – we will now extend the independent cascade model to continuous time domain.

Question 3.1 (optional). Look at the implementation of the independent cascade model in the function `GenCascadeIC(...)` of `mlss.cpp`.

Influence Estimation

Given a spreading process that started in the set of source nodes A , we define $N(A)$ as the number of nodes infected at the end of the spreading process and then define the influence function $\sigma(A)$ as the average total number of nodes, *i.e.*, $\sigma(A) = \mathbb{E}N(A)$. Now, we will simulate spreading processes and estimate influence.

Question 3.2. Generate a Forest-Fire network with 200 nodes, any forward and backward burning probabilities you like and edge weights (edge probabilities) uniformly distributed between 0 and 0.5. Then, generate 1,000 spreading processes (or cascades) over the network using a random source node set of size 2 (fixed for all cascades). To help you get started, we tell you how to generate the cascades:

```
$ ./simulate_ic -i:"your-forest-fire-network-namefile.txt" -t:1 -  
ns:5 -c:1000 -cg:25 -sg:1 -o:"forest-fire-cascades"
```

where `your-forest-fire-network-namefile.txt` should be the filename of your forest-fire network. The command will print out the (random) source set, average influence and standard error based on the simulated 1,000 cascades. It will also output several files, including two plots and a Gephi file.

Have a look at the plots, and find out what they show. The code automatically tried to fit the quantities of one of the plots to a functional relationship. Do you know which type of relationship?

Question 3.3. Open the Gephi file of the cascades you just generated in the previous question using Gephi. Then, press the bottom “Enable Timeline”, that will help you to visualize the cascades spreading. For this, look for the small shell on the left bottom and play around with different values for the custom bounds and the

play settings. Push the play button also on the top left, and you should be able to observe cascades spreading over the network for the appropriate bounds and play settings. It may help to choose the right layout.

Question 3.4. Explore the effect of choosing other source sets. In particular, try different fixed (deterministic) source sets of different sizes by modifying options `-t:` and `-ns:`. Node id's will be distributed from 0 to 127. Do you find large differences in average influence or cascade size distribution for different source sets?

Question 3.5. Generate networks with different characteristics, *i.e.*, different network types, edge densities, and edge weights and simulate spreading processes over them. Do you notice large differences in average influence or cascade size distributions? Which are the characteristics that have a greater impact?

Influence Maximization: nodes

Now that you know how to simulate spreading processes and understand them, try to find the optimal set of sources that maximize the average influence in a given network. We have implemented a baseline that chooses the K nodes with highest out-degree in `maximize_sources_ic.cpp`.

Coding 3.1. You need to think of other methods to choose influential nodes and implement them in `maximize_sources_ic.cpp`. Your code should go in the space delimited by TODO comments. You may like to rank nodes by (i) some network property you can compute directly from the network, (ii) some measure computed on simulated spreading processes, or (iii) hybrid approach combining (i) and (ii). Can you think of any related theoretical computer science NP-hard problem? This may help you to find some efficient approximate solution to the problem with provable theoretical guarantees. Your implementation should work for any number of source nodes K .

Question 3.6. Test your method(s) in one (or more) of the random networks you have generated in previous sections, generate new ones or try your method(s) in the networks located at `code/mlss-17/data/`. Note that networks in `code/mlss-17/data/` are unweighted and you will have to set an edge probability with option `-p:` both in `maximize_sources_ic` and `simulate_ic`. Once you believe you have found a good method, join the competition.

Competition. Apply your methods for finding optimal source sets of size $K = 1 \dots 10$ to the networks in `code/mlss-17/data/`. Use option `-p:0.5` both in `maximize_sources_ic` and `simulate_ic`. For every network, create 10 files, one per K , in which you write the node id's (or node id for $K = 1$) of your

solution in a single line, separated by a semicolon (;). For example, for $K = 4$ the file should contain only one line:

```
10;2;3;4
```

For the filenames, use the following format. For example, for the network `mlss-network.txt`, $K = 3$ and imagine your surnames are Gomez and Rodriguez, then the filename should be called `mlss-network-3-gomez-rodriguez`. Send the files by email to `manuelgr@tuebingen.mpg.de`.

Whenever you believe you have mastered influence maximization or you would like to try something perhaps more challenging, move on to the next section.

Influence Maximization: edges

Here, we will try to find the optimal set of edges that added to the network maximize the average influence on a network given any random source set. For the added edges, we assume they are always active, *i.e.*, its edge probability is 1.0. We have implemented a baseline that chooses K random edges starting from nodes with the smallest out-degree in `maximize_edges_ic.cpp`.

Coding 3.2. You need to think of other methods to choose optimal edges and implement them in `maximize_edges_ic.cpp`. Your code should go in the space delimited by `TODO` comments. Do you think it is an easier or a harder problem than finding the optimal source set?

References

- [1] A.-L. Barabási and R. Albert. Emergence of scaling in random networks. *Science*, 286:509–512, 1999.
- [2] D. Kempe, J. M. Kleinberg, and E. Tardos. Maximizing the spread of influence through a social network. In *KDD '03: Proceedings of the 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2003.
- [3] J. Leskovec, L. Backstrom, and J. Kleinberg. Meme-tracking and the dynamics of the news cycle. In *KDD '09: Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2009.
- [4] J. Leskovec, J. Kleinberg, and C. Faloutsos. Graph evolution: Densification and shrinking diameters. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 1(1):2, 2007.