

Social Network Analysis at Scale in Cloud

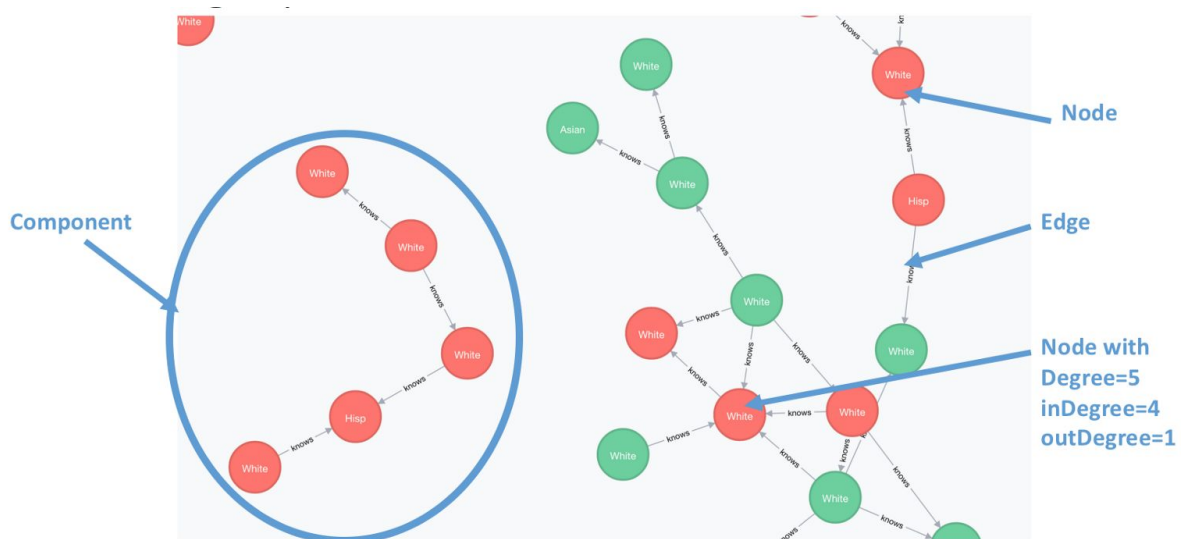
Project Proposal

Xiaolan Cai

Problem statement

Social network analysis has drawn a lot attention for the past decades. Initially sociologists have started to study how social networks are formed by and how people are connected. Nowadays social networks like Facebook, Twitter are used to spread information like news, ideas. And also for marketing purpose, for example we may also want to study how to maximize the influence of advertisement by choosing proper targets.

Social networks normally can be modeled as graph structure model, vertices/node represent an individual and edge represents the relationship between individuals.



Meaning time big scale data mining tools is well developed, distributed computing significantly improving the speed of mining data by parallel execute subtasks divided from one single big problem with distributed systems in scalable ways. And social network analysis normally

includes massive amount graph-structured data to processing, how to leverage scalable distributed computing work to handle graphed structured data is indeed a interesting topic.

In this work we are interested in a specific problem of SNA, called influence maximization [1], and we want use currently start of art distributed computing tools to exam how well we can processing the graph-structured data in terms of quality and speed.

Max influence in Social network is a interesting topic of SNA (social network analysis). In [1] this problem is formally defined in [1]. Given a Graph $N(V, E)$, and choose an ignition k number of vertices set S , what is the maximum number of influence we can achieve? In [1], two Basic Diffusion Models are considered specifically Linear Threshold and Independent Cascade models. In this work, we considering only Independent cascade models for simplicity. In this model, initially start with a set of active nodes A_0 , in each step a active nodes v will have a chance to flip his any inactive neighbor w into active one with probability P_{vw} . In the whole process node v only have a single chance to active any of his inactive node w , meaning he cannot make any further attempts no matter he is successful in flip w or not. And the process stops when there is no mare activation are possible [1].

Finding the optimal solution is NP-hard problem, and in [1], an approximation solution is proposed, it use greedy algorithm to achieve slightly better than 63% of the optimal solution. However, the computation is still pretty heavy, because they need to run Monte Carlo simulation in each update step?

It is challenging task to processing massive data not even to say graph structured data. There are two kind of graph processing framework. GraphLab: share state in parallel vertex computation. Pregel: passage passing, graphx

We need modified existing algorithm to make it suitable for vertex programing, which is request by most graph based commutating framework. We provide a vertex program that is run by each vertex in parallel.

So our question is how influence maximization problem can be processed in cloud computing environment with scalable setups? What performance need to be sacrificed in terms of accuracy to archive speed? (what's the trade off between quality and speed?)

Solutions for processing

One option is to use spark graphx which use message passing between vertices (e.g. Pregel). We first store our data in HDFS and use graphx to build graph by create vertices and edge from

it. As normal spark programming use. map to process the source data and then doing targeted programming computing on individual graph objects.

Another option is to combine spark with graph based database such as Neo4j, we import the dataset into Neo4j first. Then we need distribute the data in neo4j to HDFS, after it spark will do parallel computing and put result in HDFS and followed by exporting back to neo4j.

Beside message passing based graph computing framework, we can also use stat shared tool, such as graphlab, we can also explore this tool to for comparing.

Another orthogonal method of improvement the speed of computing is edge sampling, for example we only delete some of the edge of a vertex that have high degrees to reduce the computation with approximation [3].

Setups and data sets:

We will setup our experiment in Cloudlab.

Dataset, we will use undirected social network in snap datasets, such as “ego-facebook”, “gemsec-facebook” [2]. And the data comprise lists of pair of node-ids, which represents edges of graph. The data file is in txt or csf format. We will store the data in HDFS in our cloudlab setups.

Timelines of the project:

We can divided our project in following phases

Phase 1 (1 week)

We will setup the experiment in AWS/Cloudlab for 1 week, including spark graphx alone setup and neo4j+spark setup.

Phase 2 (1.5 weeks)

We will measure spark alone setup

Phase 3 (1.5 weeks)

Measure neo4j + spark setup

Phase 4 (1 week)

Report write up.

Challenges:

To adapt the influence maximization problem for spark like framework is difficult, as we need to define vertex based programming for parallel with approximation and not sacrifice the accuracy too much.

Issues:

Time constraint, since I am doing this project only, I have concern about I may not be able to finish all the measurement I have planned to, but at least two feasible scalable setups should be built and measured.

Spark GraphX uses Scala for programming and I need to learn Scala for this.

Reference:

1. Maximizing the Spread of Influence through a Social Network
2. <https://snap.stanford.edu/data/>
3. Mark S. Handcock, David R. Hunter, Carter T. Butts, Steven M. Goodreau, and Martina Morris (2003). statnet: Software tools for the Statistical Modeling of Network Data. URL <http://statnetproject.org>
4. Bridging the GAP: Towards Approximate Graph Analytics