

## ##

This is just a summary for a very rough and preliminary study about network detection theory advised by Prof. Santo. Raw data and block number was collected/calculated by Lucas.

### **Test john55 dataset.**

**John55.gml is a social network with about 5000 nodes and many edges.**

#### **Introduction:**

Stochastic block model is a method to detect communities in a network.

Find the necessary information for stochastic block model in this paper:

<https://arxiv.org/abs/1008.3926>

We are going to test if the stochastic block model can work for most practical data in real world. Here is what we want to test and the procedure:

1. Find 2 communities A and B detected by stochastic block model.
2. For community A, find 2 nodes i and j. Calculate  $K_i$  and  $K_j$ , where K is number of the nodes connected to them (the degree of a vortex).
3. Find  $K_{i_B}$  and  $K_{j_B}$ , where  $K_B$  is the number of neighbors in community B.
4. Plot  $K_{i_B}/K_{j_B} \sim K_i/K_j$ , if it's a straight  $y \sim x$  line?

#### **First, cluster the network**

##### **1. Clustered by `graph_tool.minimized_blockmodel_ml()`**

Instead of my naïve implementation of Newman's stochastic block model, we can just use igraph. I'll use 2 igraph methods and then show the  $K_{i_B}/K_{j_B} \sim K_i/K_j$  plots for each.

graph-tool python package is fast with help of some C/C++ library so the method is almost linear  $O(N \log N \log N)$  and it's the only available stochastic block model (SBM) method I can use for now instead of the code I wrote for SBM.

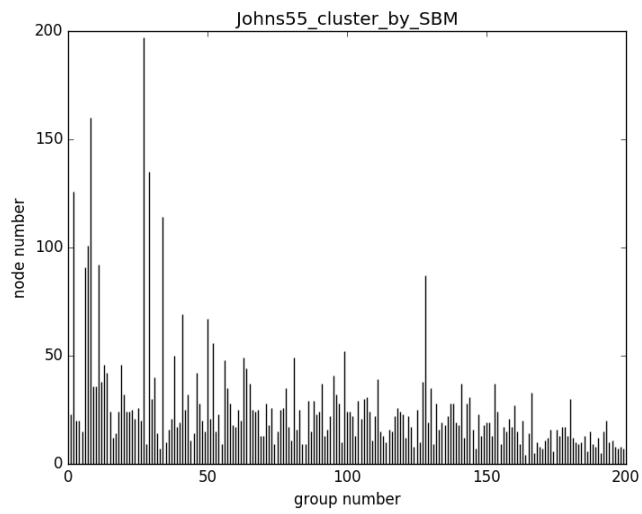
Document for the function.

<https://graph-tool.skewed.de/static/doc/demos/inference/inference.html>

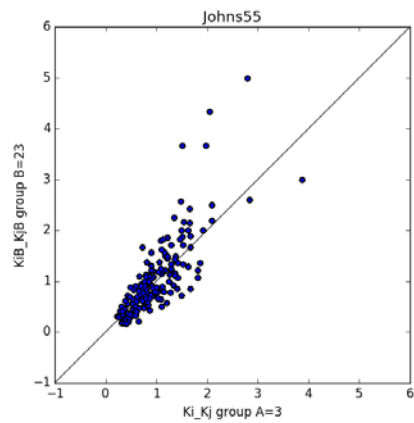
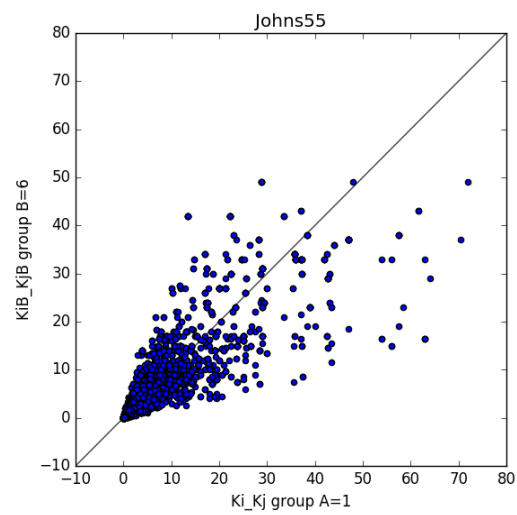
Idea is the same, maximizing likelihood  $P(G|\{b_i\})$ ,  $\{b_i\}$  is the set of parameters for block, but the procedure is different from Newman's 2011 paper. It doesn't always pick the move with largest  $\Delta L$ . It uses the method in this paper:

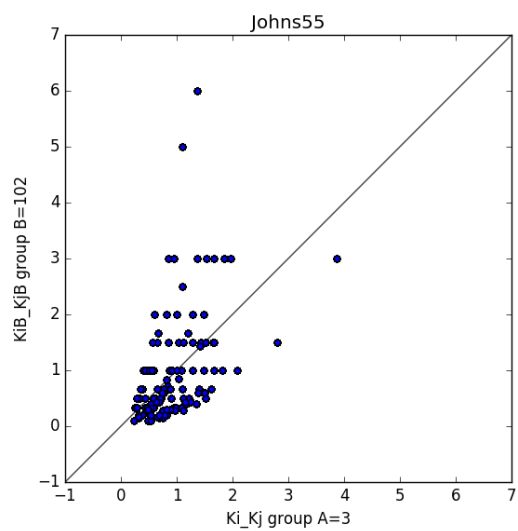
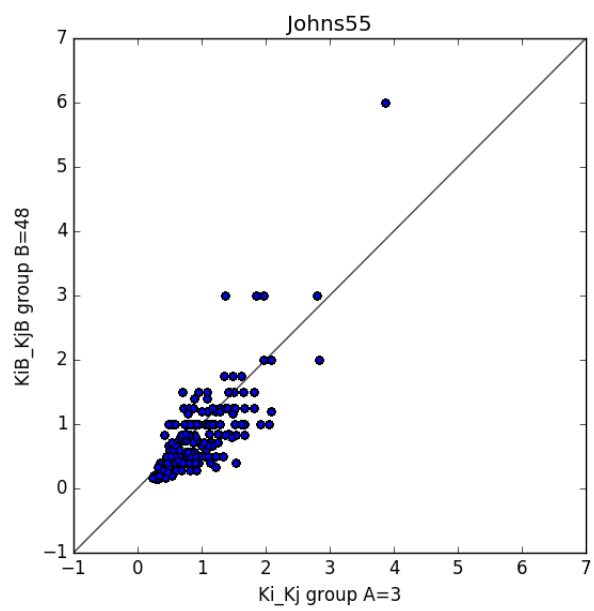
Tiago P. Peixoto, "Efficient Monte Carlo and greedy heuristic for the inference of stochastic block models", Phys. Rev. E 89, 012804 (2014), [DOI: 10.1103/PhysRevE.89.012804](#), [arXiv: 1310.4378](#)

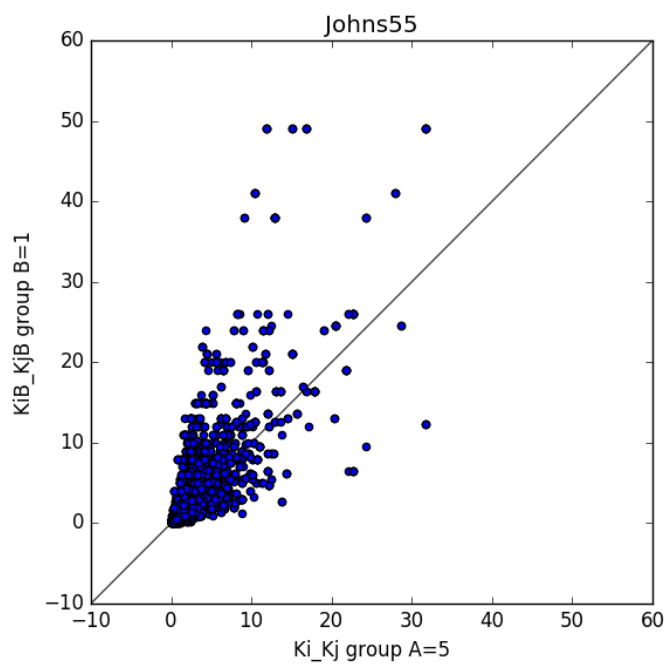
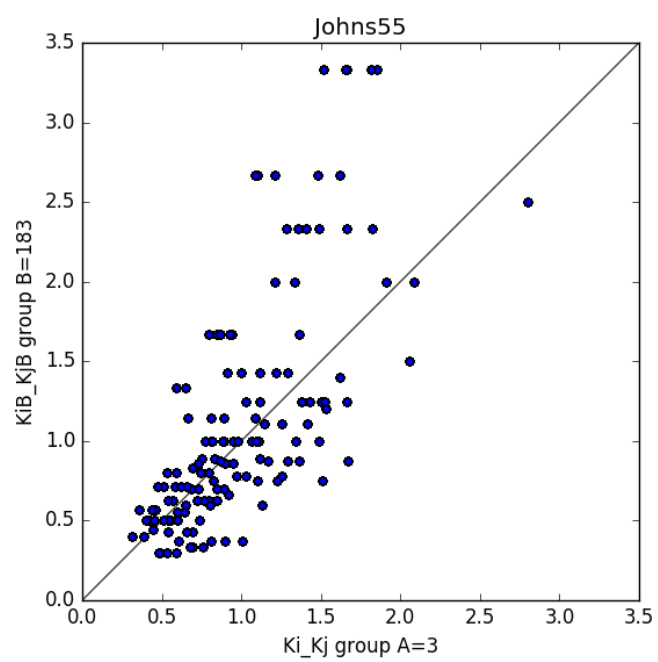
The method was tested in network with size from  $N=369$  to  $N=654$  782 in the paper.

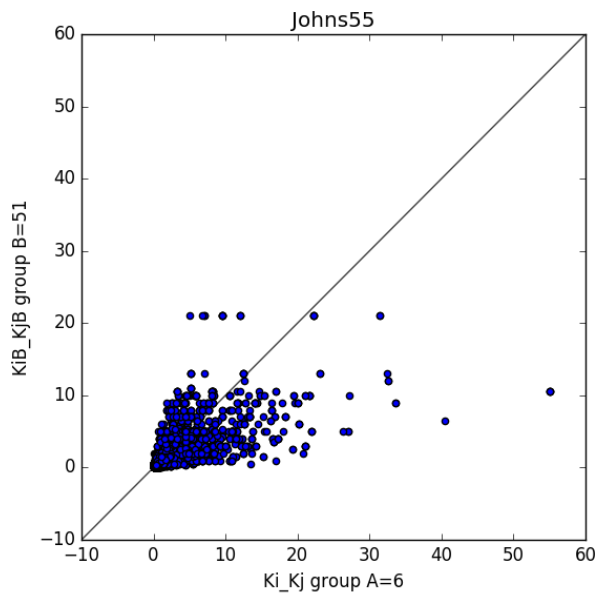
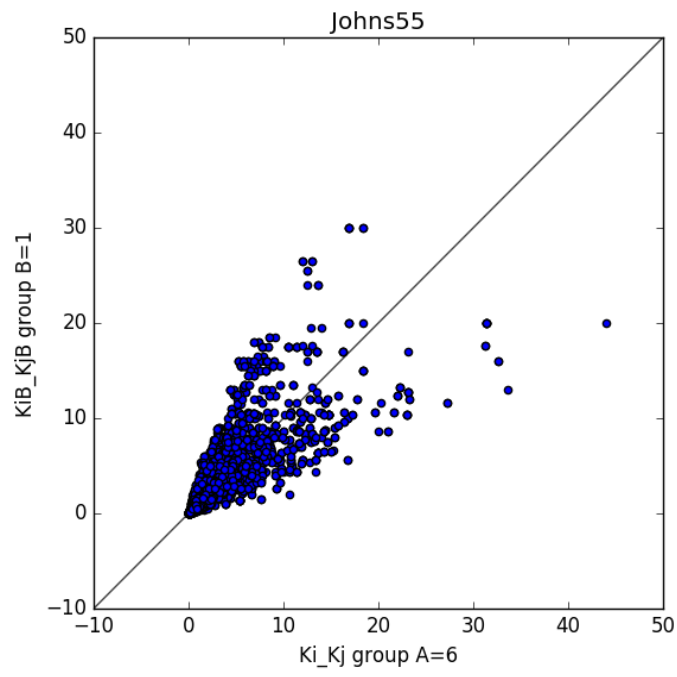


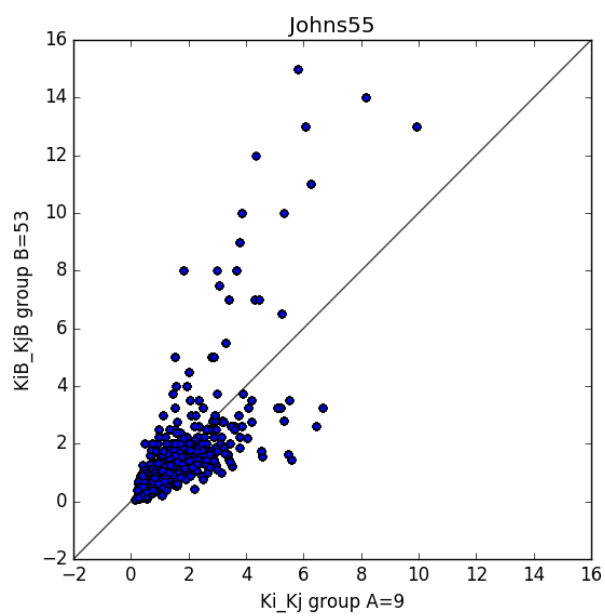
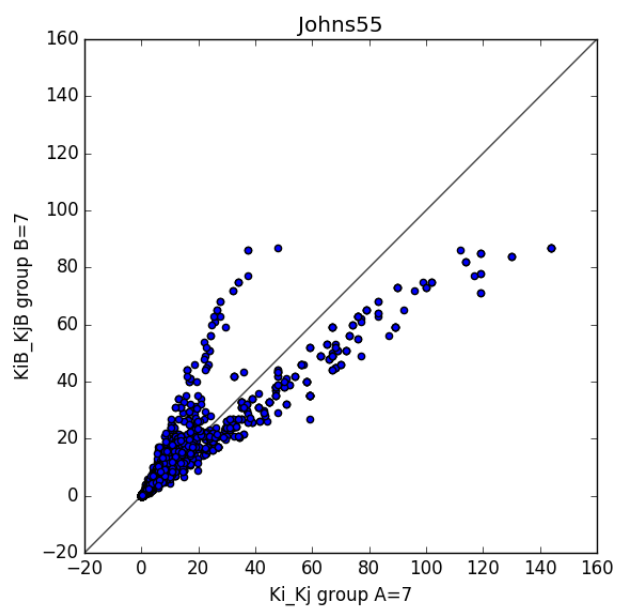
**Then plot the  $K_{iB}/K_{jB} \sim K_i/K_j$**   
list some significant plots:

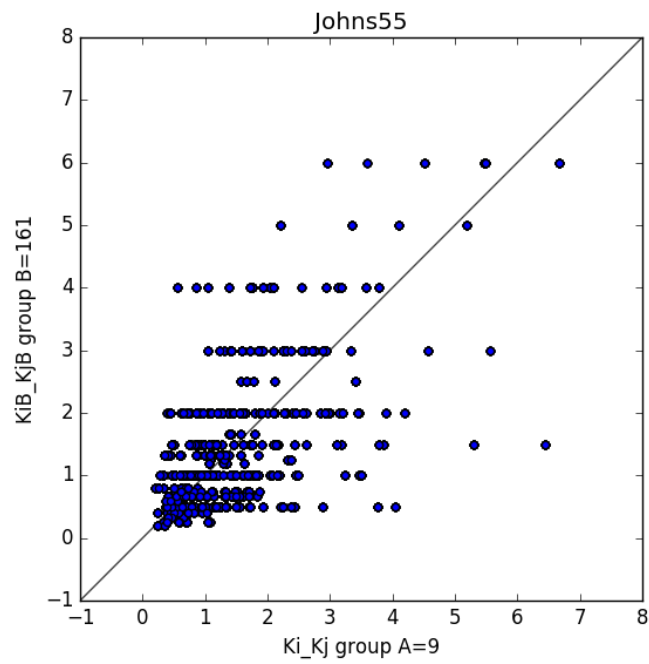
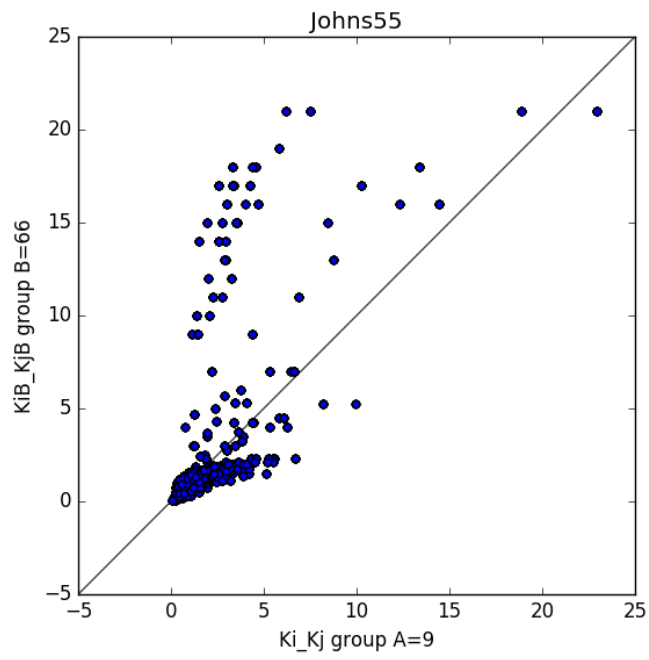


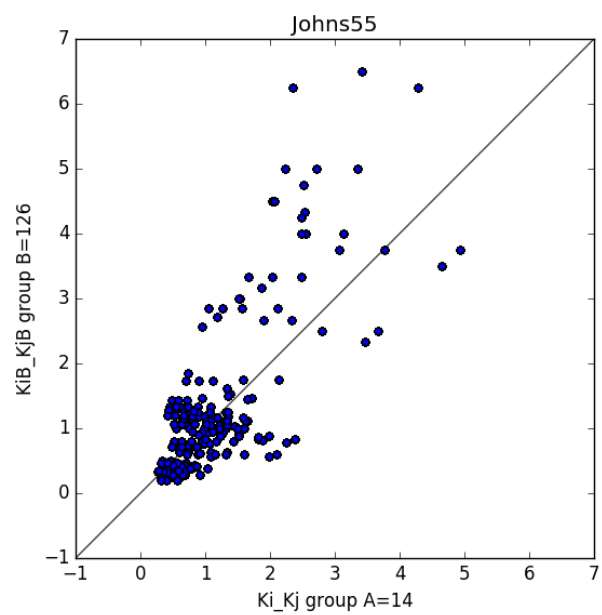
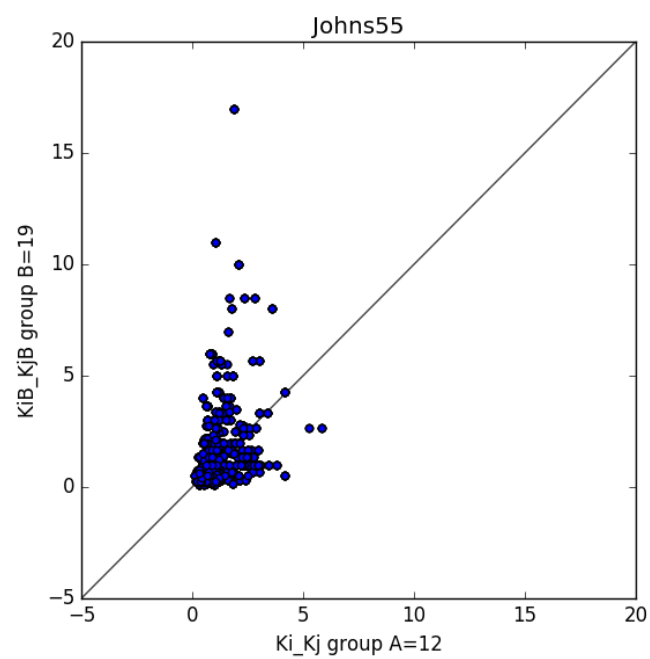




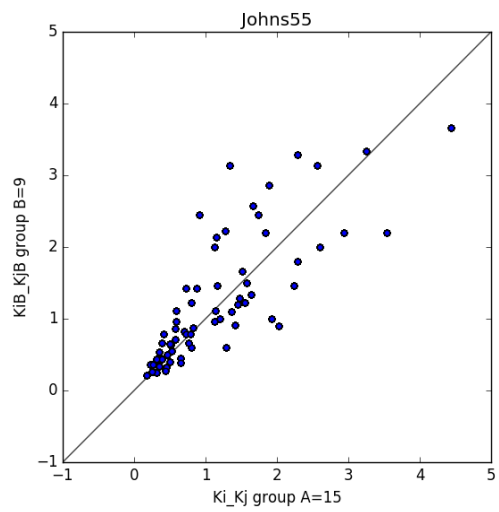
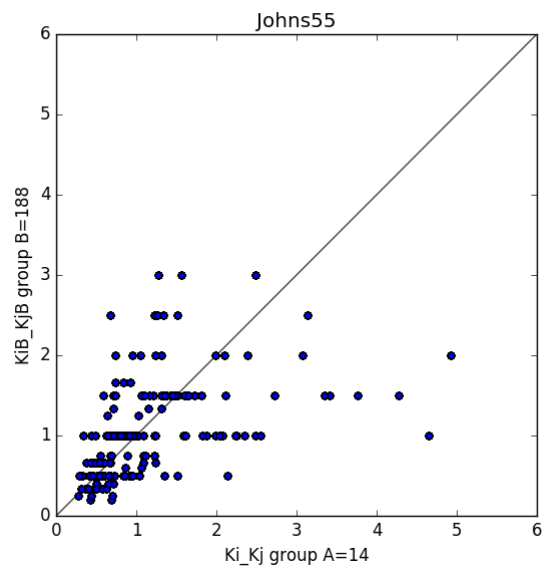


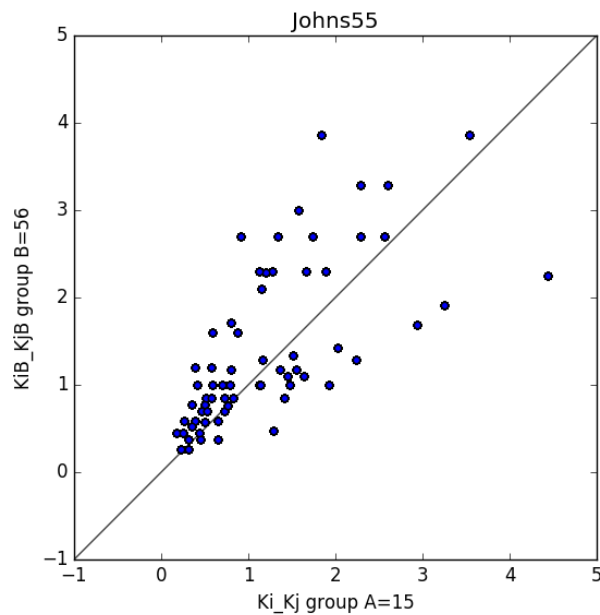








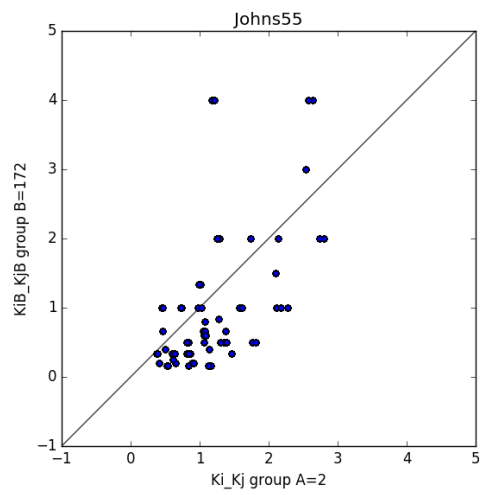
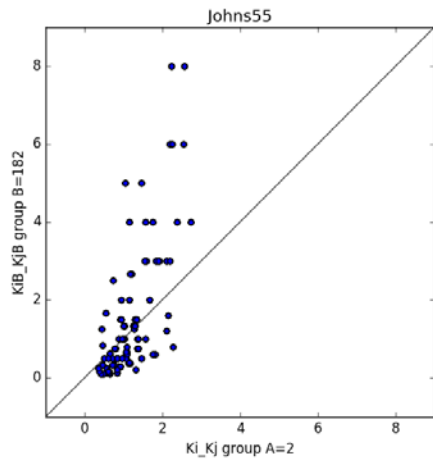
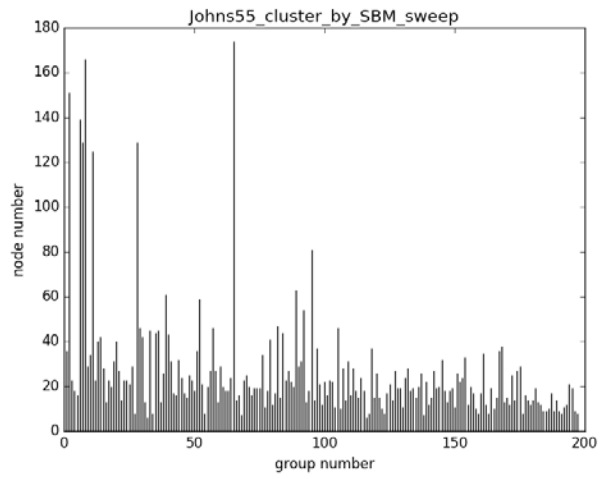


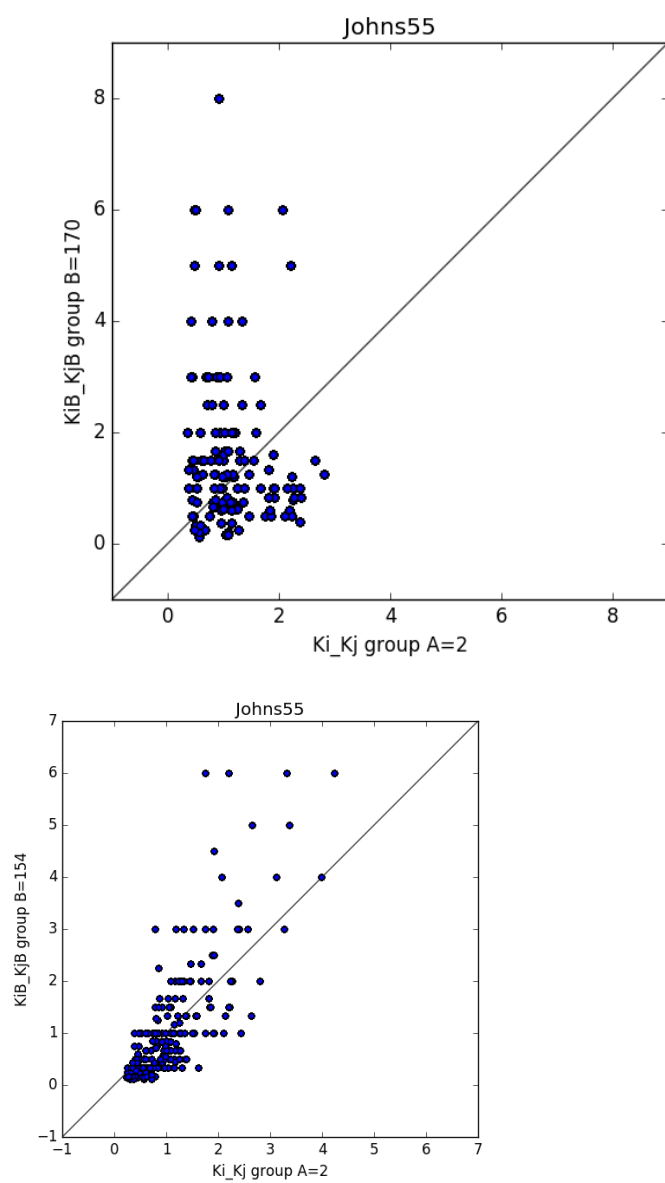


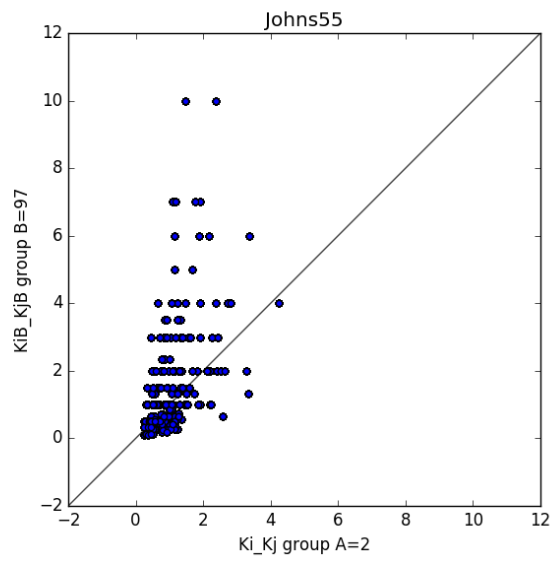
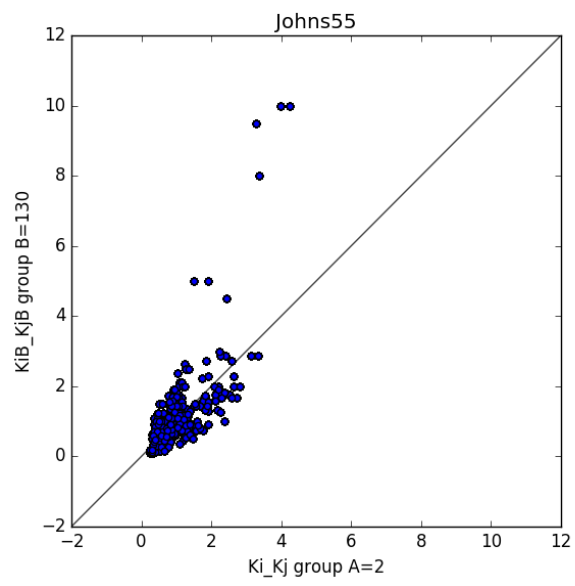
## 2. Clustered by `graph_tool.minimized_blockmodel_ml()` and `mcmc_equilibrate()`

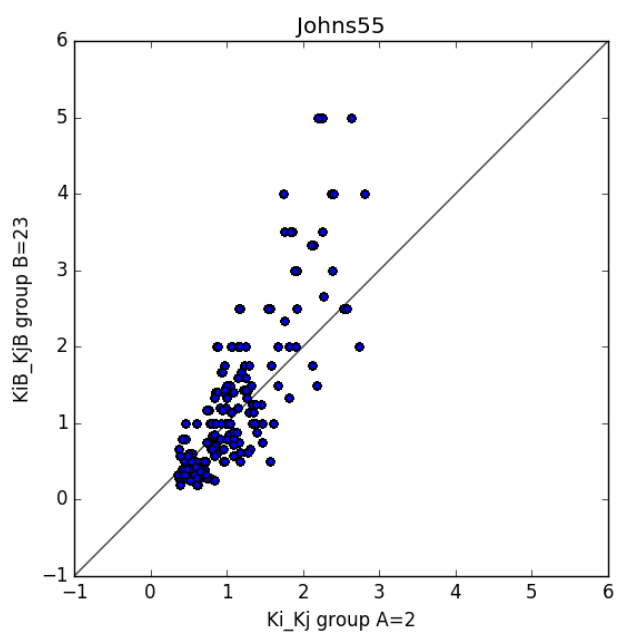
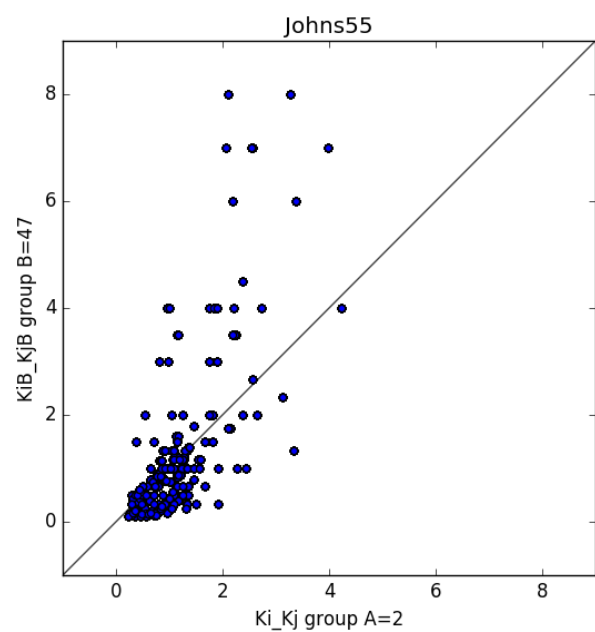
However `graph_tool.minimized_blockmodel_ml()` cannot correctly cluster `karate.gml` into 2 groups, it always results in one single group. Maybe it thinks the whole thing is more likely to be a single community, which makes me suspicious on its correctness.

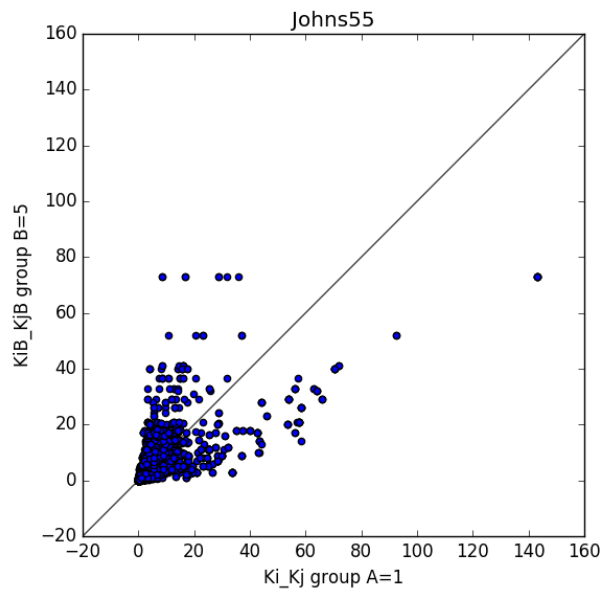
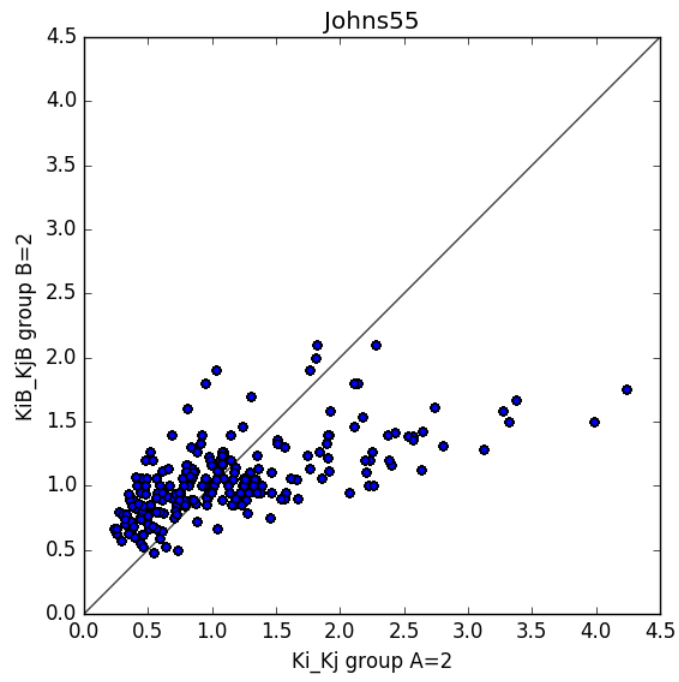
An improvement is `mcmc_equilibrate()`, which can avoid metastate, and use the state obtained by `graph_tool.minimized_blockmodel_ml()` as initial state.

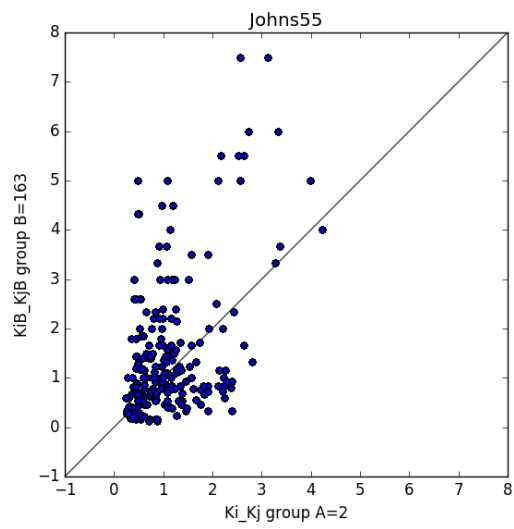
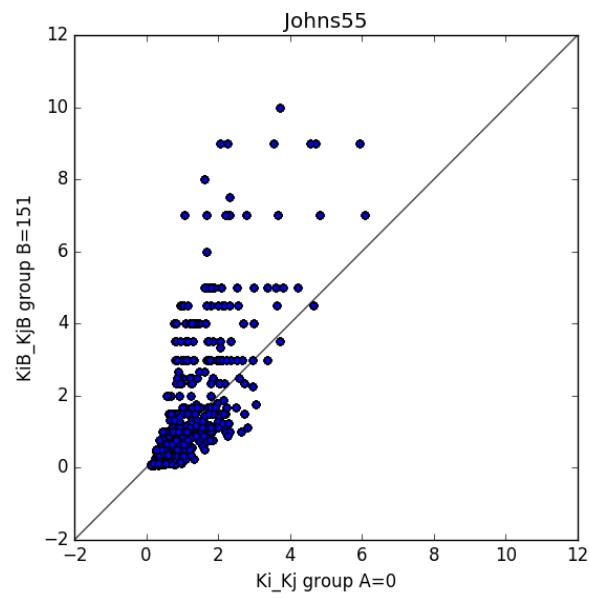




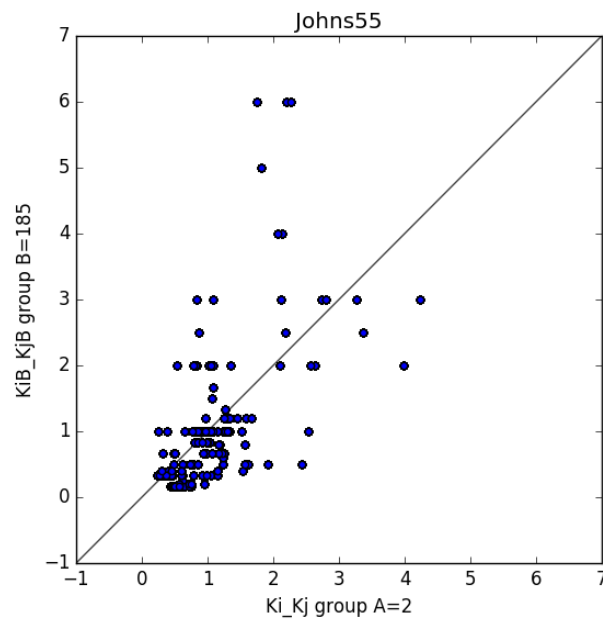








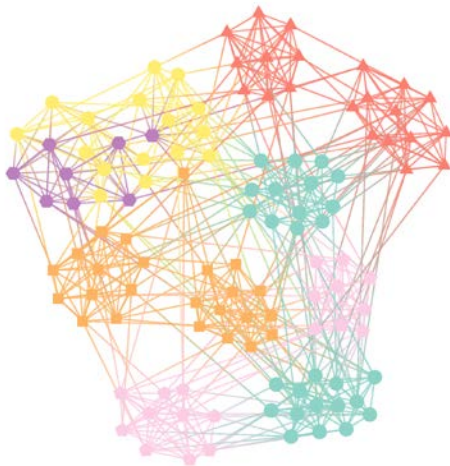




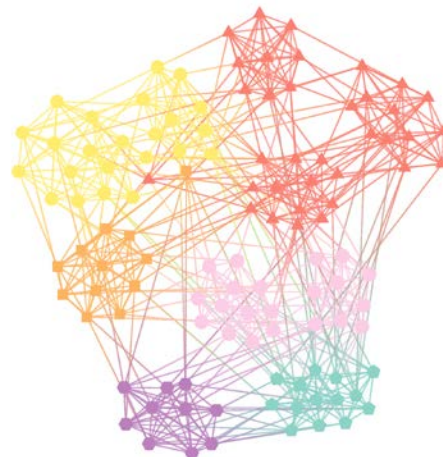
### Test both methods using karate.gml and football.gml

Neither can cluster karate.gml correctly. Maybe it's because the size of network is too small.

The second method is better for football.gml. Divided it into 6 groups.



first method



second method