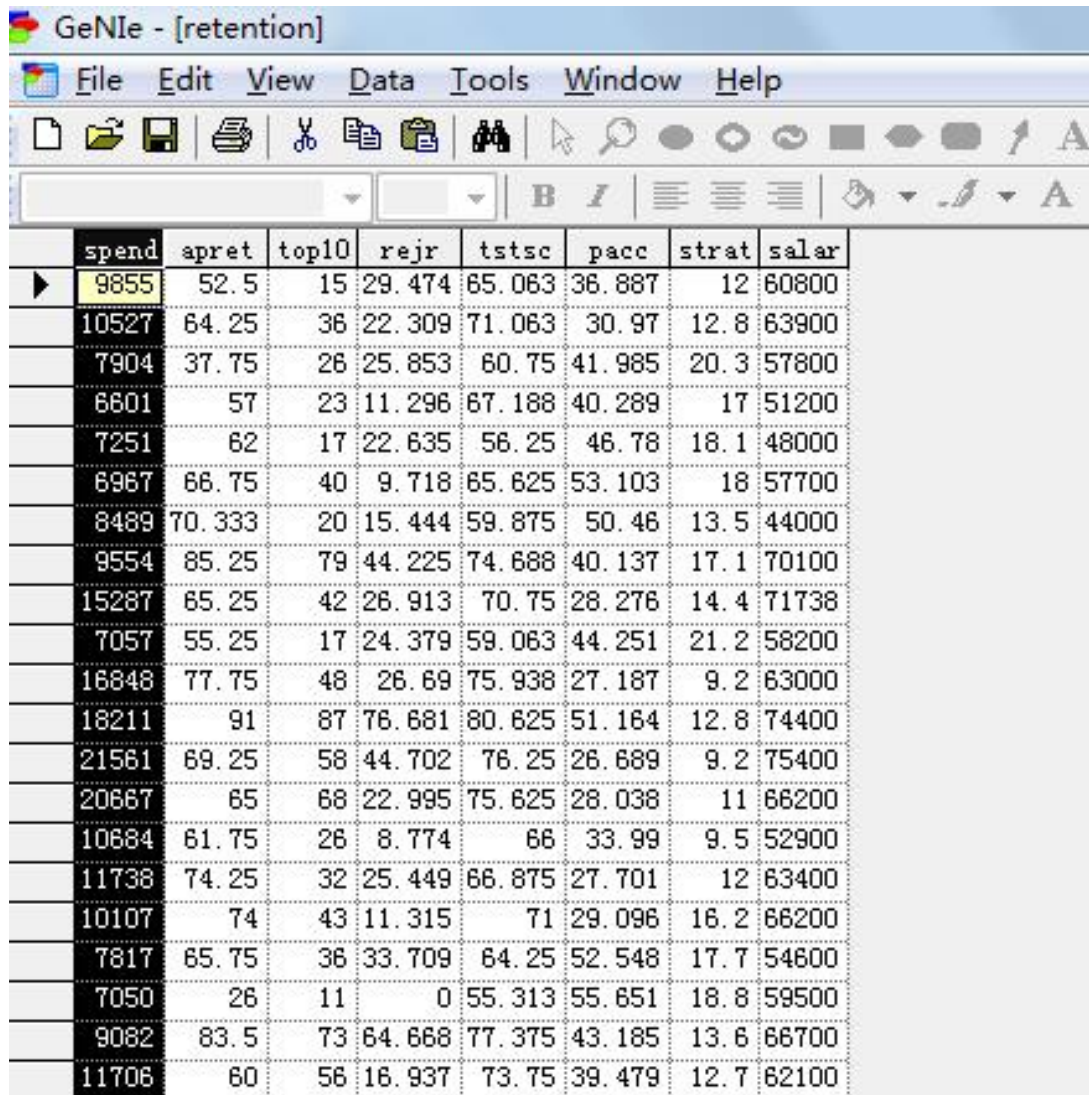# Assignment 5

## Yue Su (yus55)        Dan Sun (das225)

### 1. Load the data

Load the data into system. Screen shots below shows the process.
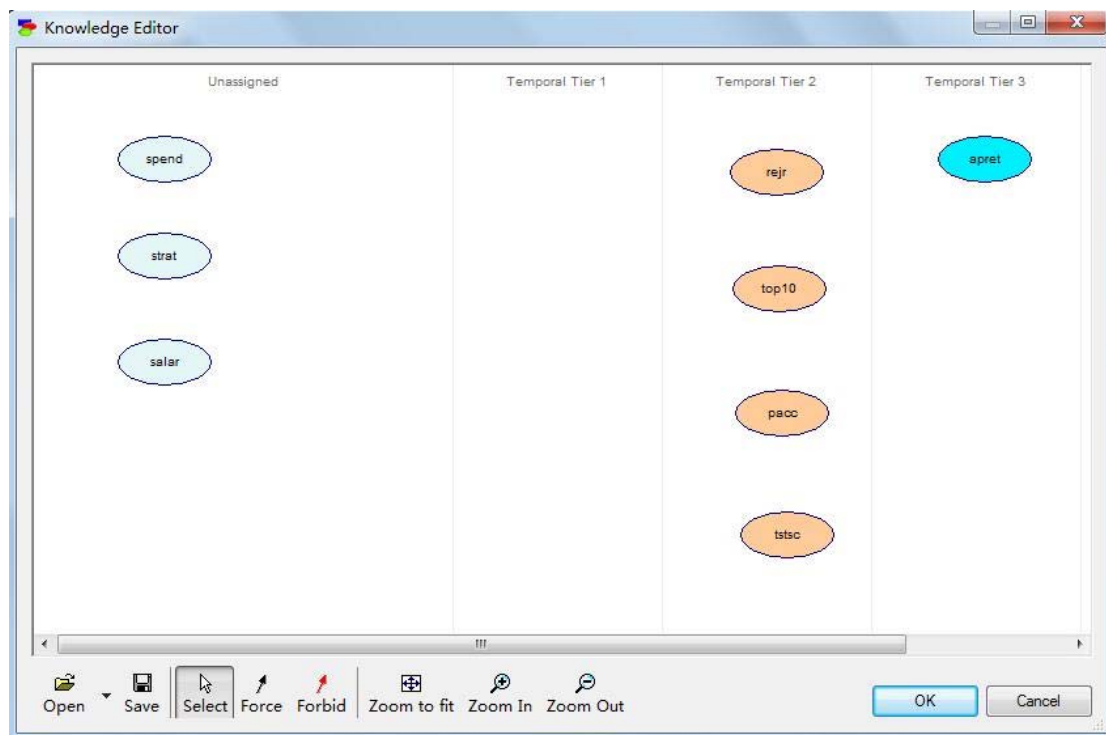


**2. Enter the background knowledge by clicking the Background Knowledge button. Here you can force and forbid causal connections and also order variables in temporal tiers.**
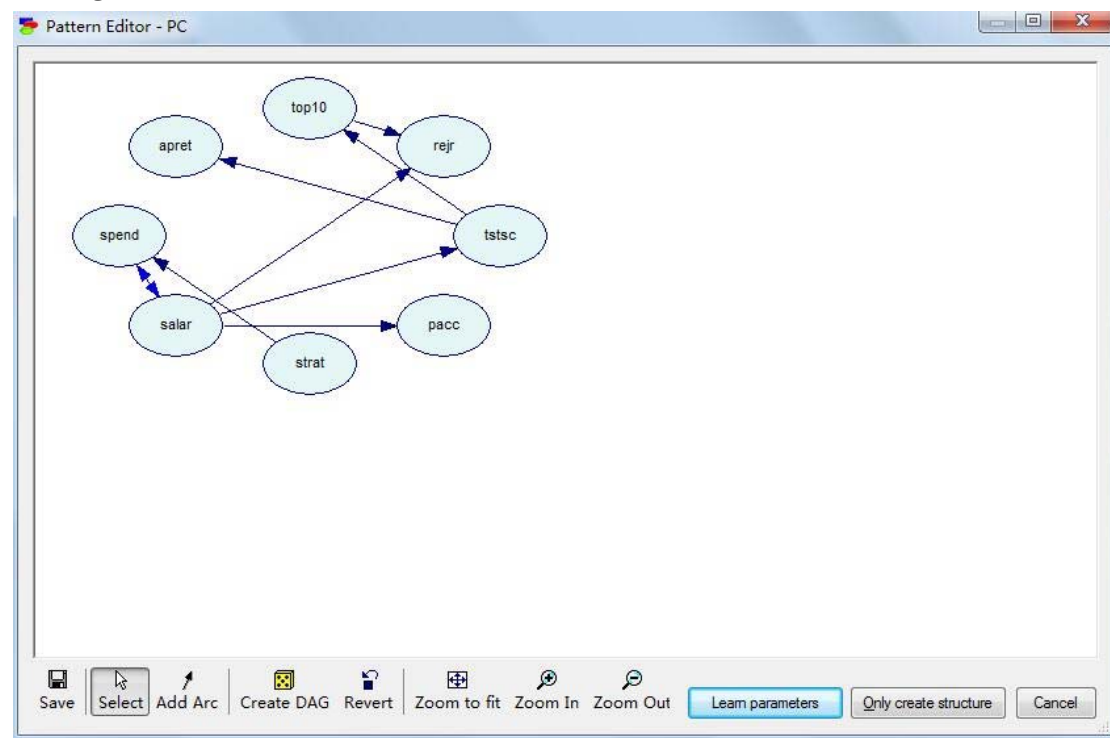
**3. Choose the PC algorithm and set the significance level (PC algorithm uses classical statistical independence tests and here is where you set the significance level for these tests)**
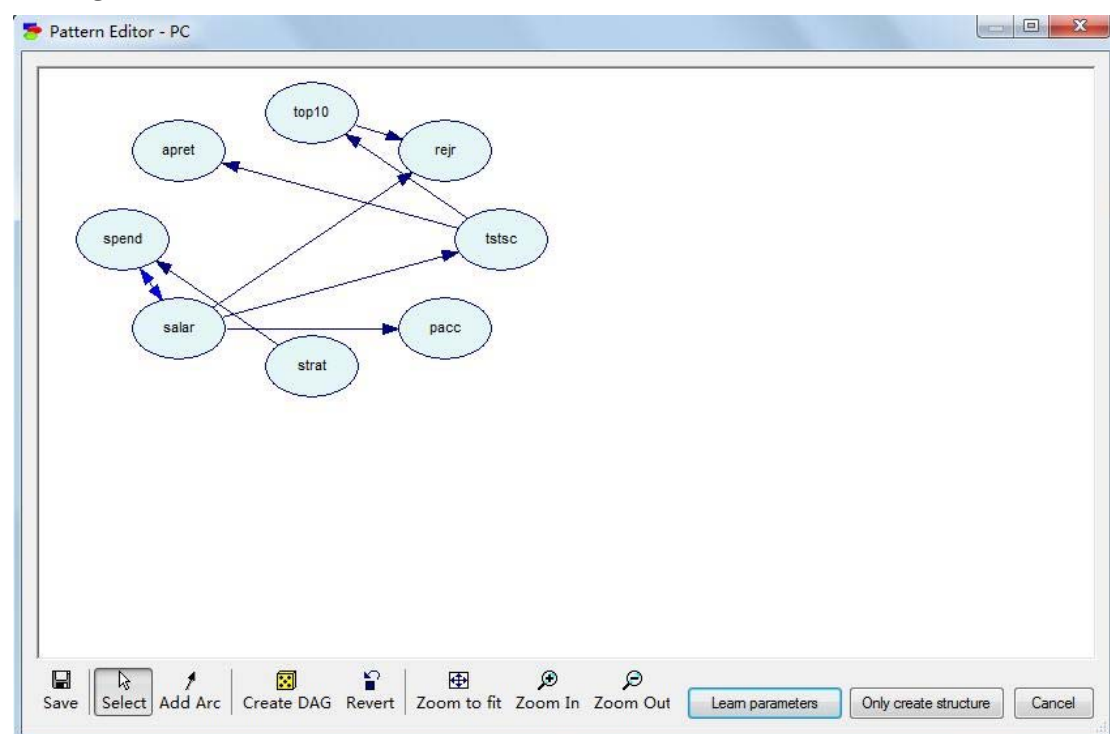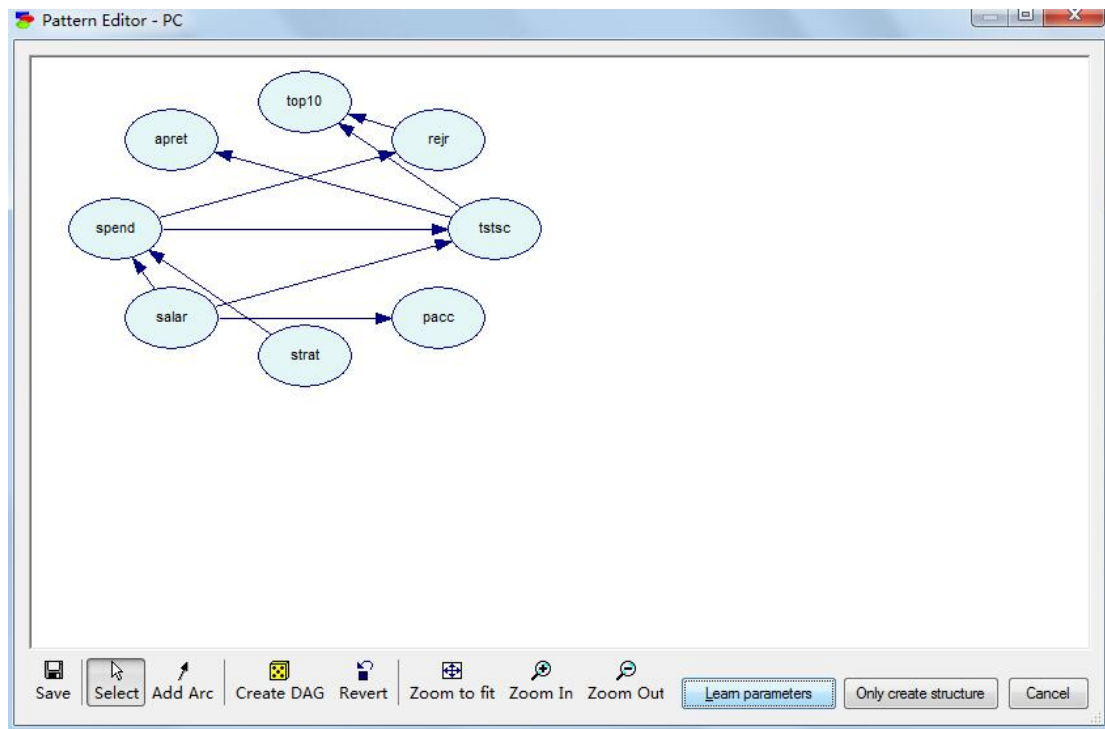
- Significance level = 0.0001

- Significance level = 0.0005



- Significance level = 0.001

- Significance level = 0.005



Depending on the significance level used in tests above, statistical decisions regarding independence may be different and a different class of causal structures may result. It is, therefore, a good practice to run the program at several significance levels. We ran algorithm with the following significance levels: p = 0.0001, 0.0005, 0.001, 0.005. The core of the structure, i.e., how freshmen retention rate and graduation rate are related to the remaining variables, was insensitive to changes in significance. This suggests that the structure proposed is robust. The edges of the graph have the following meaning: A single arrow (⟶) denotes a direct causal influence. A double headed arrow (⟷) between two variables denotes presence of a latent common cause of these two variables. An single arrow with a circle at one end (o⟶) expresses inability to deduce whether there is a direct influence between the two variables (⟶) or a latent common cause between them (⟷). An edge with circles at both ends (o�—o) expresses inability to deduce whether there is a direct influence between the two variables and, if so, what is its direction or a latent common cause between them (⟷).

Most graphs contained a direct causal connection between the average test scores and freshmen retention. Also, the graphs contain latent common cause connection between spend and salary.