# Assignment 5 Causal Discovery

HAITAO  ZHOU  (HAZ59)

YINGZHI  YANG(YIY50)

XIN  JIN(XIJ21)

In this assignment, we use GeNIe to analyze the data.

1. We import the data retention.txt into GeNIe:

2. We can learn from some statistics of these natural data:

| | Mean | Variance | StdDev | Min | Max | Count |
|---|---|---|---|---|---|---|
| spend | 10974.5 | 3.02507e+07 | 5500.07 | 4125 | 35863 | 170 |
| apret | 56.7211 | 326.781 | 18.0771 | 18.75 | 95.25 | 170 |
| top10 | 38.4588 | 547.859 | 23.4064 | 8 | 98 | 170 |
| rejr | 30.6542 | 292.345 | 17.0981 | 0 | 84.067 | 170 |
| tstsc | 66.1642 | 48.6549 | 6.97531 | 48.125 | 87.5 | 170 |
| pacc | 43.1731 | 171.746 | 13.1052 | 8.964 | 76.253 | 170 |
| strat | 16.0865 | 16.0521 | 4.0065 | 7.2 | 29.2 | 170 |
| salar | 61357.6 | 9.60946e+07 | 9802.79 | 38640 | 87900 | 170 |

3.With the correlation matrix, we can see the correlation ratio between every two features. We can tell from this chart that apret may be positively correlated to top10, rejr, tstsc, salar and negatively correlated to pacc and strat.

| | spend | apret | top10 | rejr | tstsc | pacc | strat | salar |
|---|---|---|---|---|---|---|---|---|
| spend | | | | | | | | |
| apret | 0.601231 | - | | | | | | |
| top10 | 0.675656 | 0.642464 | - | | | | | |
| rejr | 0.633544 | 0.514958 | 0.643163 | - | | | | |
| tstsc | 0.71491 | 0.782183 | 0.798807 | 0.628603 | - | | | |
| pacc | -0.283673 | -0.102834 | -0.207505 | -0.0715207 | -0.164223 | - | | |
| strat | -0.581755 | -0.458311 | -0.247857 | -0.283617 | -0.465226 | 0.131858 | - | |
| salar | 0.711838 | 0.635852 | 0.637648 | 0.606777 | 0.715472 | -0.37524 | -0.347673 | - |

4. To have an overview of the data, we develop the pattern without setting any background or changing the significant level. In this case,

`apret=1.88181*tstsc-0.543687*strat+Normal(-59.0413,11.0966)`



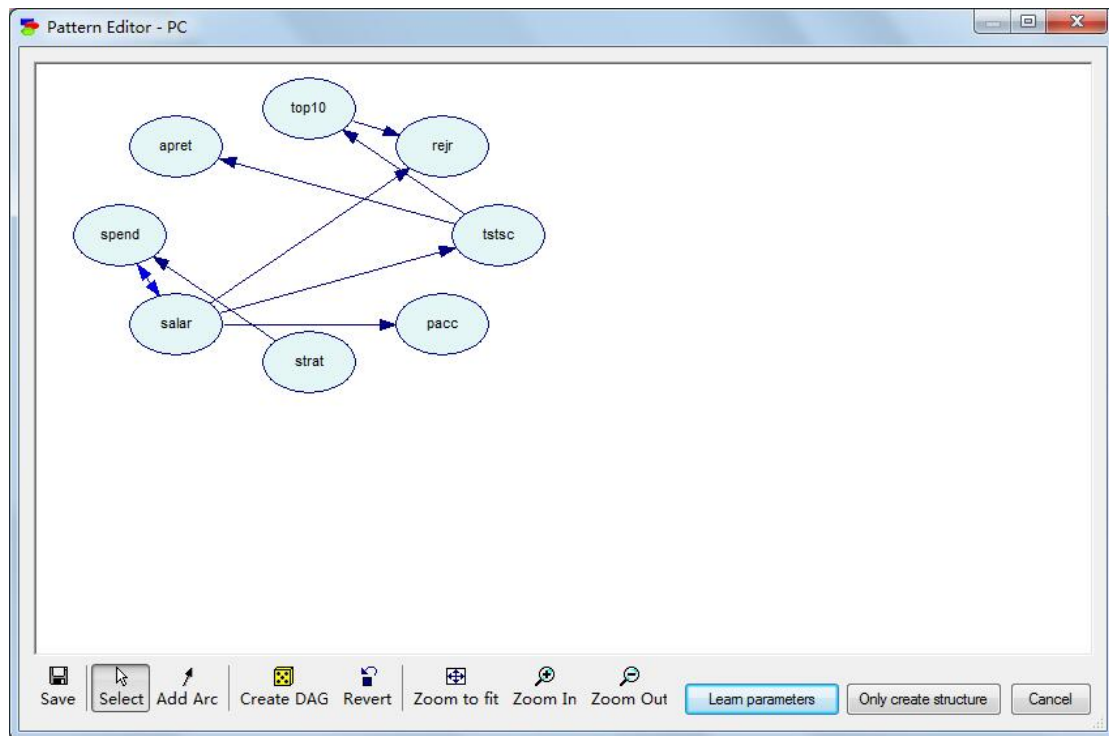5. In order to force and forbid causal connections, we separate these 8 features into 3 temporal ties, which is "spend, start, salar" , "rejr, pacc, top10, tstsc" and "apret".

6.Furthermore, we train the data by setting different significance level into pattern editor, and watching the connection's change:

(1)significance level = 0.001:
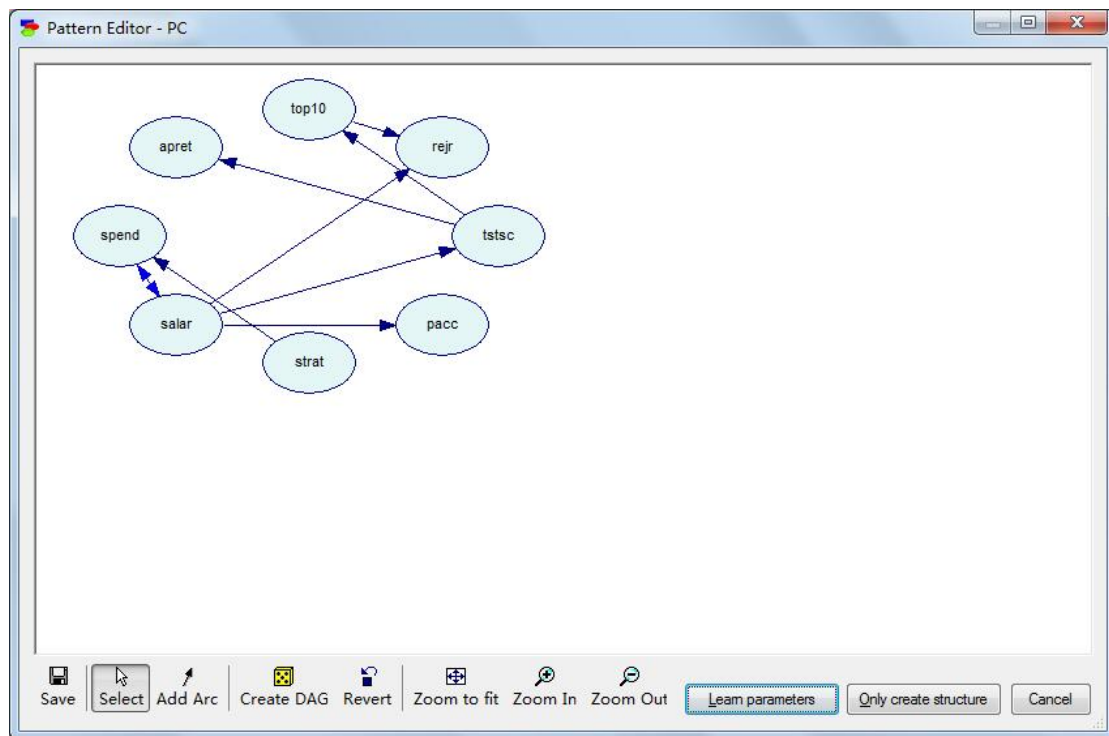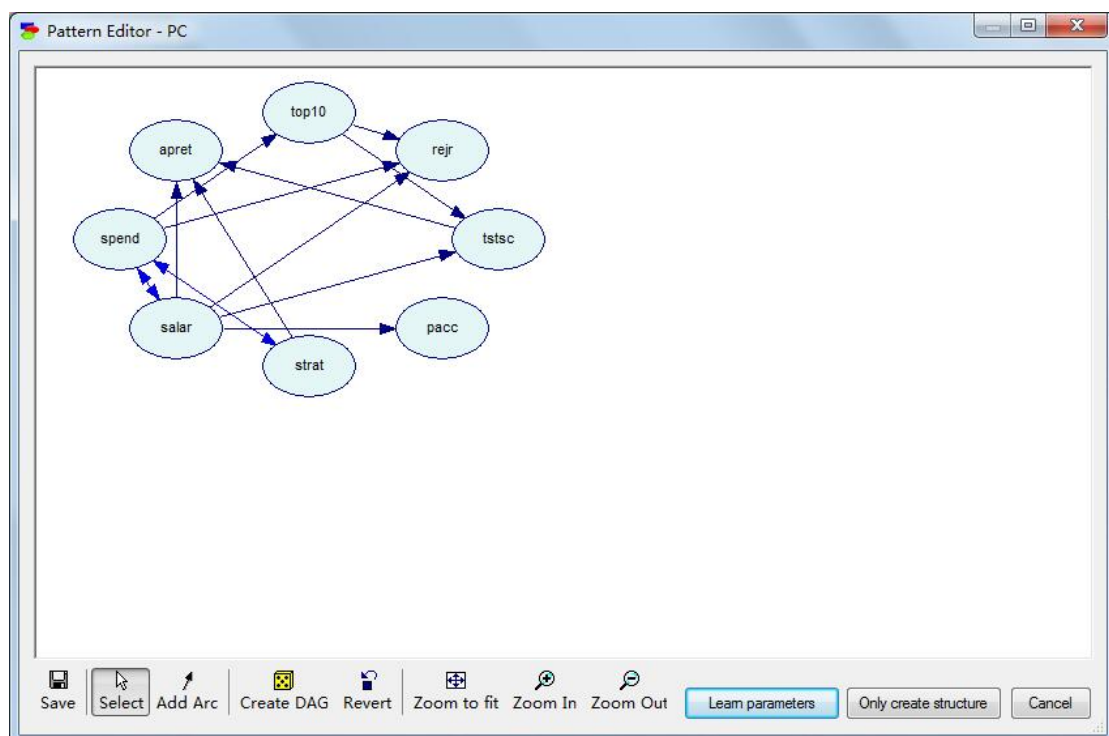
`apret=2.02709*tstsc+Normal(-77.3999,11.2629)`



(2)significance level = 0.01:

`apret=2.02709*tstsc+Normal(-77.3999,11.2629)`

(3)significance level = 0.1:

apret=1.60236*tstsc-0.530669*strat+0.000281388*salar+Normal(-58.0262,10.9281
)

Conclusion:

The diagrams showed that "tstsc" affects the "apret" mostly. In Druzdzel & Glymour's conclusion, "top 10" and "tstsc" are the factors determine the value of "apret". And in our research, when the significance level is low, the "top 10" factor will not affect "apret" directly, but it does affect the "tstsc", and "tstsc" always affects "apret". And when the significance level becomes as high as 0.1, "strat" becomes another important factor which will influence "apret" and "salar" has very little effect on "apret". So, as far as we concern, the test score is the main reason for retention, and the quality of high school education will affect the freshman year test score for the college student. On the other hand, we can also see that the lower the significance level is, the more accurate the relationships are. To sum up, we believe that the data support the Druzdzel & Glymour's conclusions in some way.