

How Education Quality Affects School-Related Crime

Machine Learning for Cities

Baohan Liu, bl3360@nyu.edu; Fanshu Li, fl2301@nyu.edu;

Jacob Jiang, xj2128@nyu.edu; Wangtianhan Pang, wp2120@nyu.edu

May 5th, 2022

Abstract

The NYPD is making an average of around four arrests and seven summons daily in New York City public schools (Grace Chen, 2020), which the schools called “Pipeline To Prison”. That this situation varies largely in different schools raises questions about what attributes of schools have a great impact on the number of crimes that occur in school. We used clustering, Random Forest, and Bayesian Network to explore the relationship structure between crime in schools and attributes of schools and make predictions for crimes in school with data from the NYC school quality report and the school safety report. The result shows the school type, number of chronically absent students, and financial situation of students determine mostly whether there are crimes in school. Our result could help focus on the affective aspects to improve the management of schools and provide a safe education environment.

Background

There were 1,171 arrests and nearly 900 summonses in schools in the 17-18 school year (NYCLU, 2019), sending hundreds of young people into the criminal justice system for minor misbehavior, in which around 95 percent of the arrested were black or Hispanic, despite that there are only around 71 percent of the students are colored races (Grace Chen, 2020). And according to Student Safety Act Reporting (NYCLU, 2019), schools are used by the NYPD as a place to locate and arrest young people, which makes the kids in trouble stay away from school. Thus, the Student Safety Act (SSA) requires the NYPD to publish quarter reports on arrests, summonses, and other police-related incidents in NYC Public schools (NYCLU, 2018).

Schools can be a particularly important environment in which criminal activity occurs. Schools in high-poverty neighborhoods often struggle to retain effective teachers and have high rates of violence and dropout (Lankford, Loeb, and Wyckoff 2002; Cook, Gottfredson, and Na 2010). Also, School quality explains more of the impact on high school students.(Deming, 2011). Based on these literatures and since the crimes happen in schools, we make the hypothesis that there may be a relationship between the crimes and the environment in which the crime took place. Therefore, We decide to find out how the education quality and other attributes of schools affect students' crime.

Problem Definition

There are two questions: What attributes of schools are closely connected and interact with each other? Which group of attributes could best be used to predict whether there are crimes occurring in school? Accordingly, we have two main goals. The first is to figure out the inherent relationship structure of all the data variables, which would be achieved by clustering through different groups of variables, learning the Bayesian network from the data, and drawing an important index heatmap with Random Forest. The causal chains would help improve the whole system including the school education and management and the arrest behaviors in school. The second goal is to make as accurate predictions as possible for crimes in school through selecting appropriate models, selecting the most affective variables, building specific models for

sub-dataset, and adjusting the models by tuning hyperparameters and avoiding overfitting. The accurate prediction of potential crime in school could send an alarm for the school operation, thus the hidden trouble may be solved by proper adjustment strategies before it happens.

Since our dataset has over 38 variables and about 1300 records, clustering, Random Forest and Bayesian networks would be great tools to explore, summarize, and deal with that. In addition, high interpretability of Random Forest and Bayesian network, which also have great accuracy, would also help us understand the big picture.

Finally, the prediction models and different algorithms would be compared to evaluate the results.

Data and Preprocessing

We collected the 2017-2018 School NYPD Crime datasets from the NYC OpenData platform. The dataset includes crime records for 1919 different schools with a unique DBN in NYC. We selected and calculated the major, other, none, property, and violent crime numbers and their sum from each school in the dataset as the dependent variables.

We found the 2017-2018 School Quality Reports for Elem, Middle, High, and K-8 schools. Since the Quality rating standard for the High School is a little different from other types of schools, we selected 30 common standards for all types of schools as the Features. Meanwhile, we replaced all the string type data in the Quality report with the integer labels and used `get_dummi` to switch School Type into four different features. And we merged the crime data and school quality data on DBN and got our final dataset with 1351 data, six dependent variables, and 33 quality standards as independent variables.

Data Analysis

Firstly, we analyzed the distribution of the number of different types of crimes. According to the distribution histograms of Major, Other, No Crime, Property, Violent, and Total crime numbers, we found that all the crime types have a similar frequency distribution. Meanwhile, we found schools without any criminal records take about 75% of the dataset in all crime types. Therefore, we decide not to remove the outliers in the number of crimes. Instead, in the number of crimes, we will replace 0 with 'No crime record' and other values with 'With crime records' for further modeling.

Since the distributions are similar between different types, we chose to pick the total number as the representative dependent variable for further analysis. Then we examine the correlations among our 34 features. According to the correlation heatmap, some of the features have similar correlations and can be classified as a larger category.

Cluster Observation

We first completed the data preprocessing part. We selected the 2017-2018 dataset with a relatively complete amount of data, and filtered out the school type as high school, then merged NYPD Crime Report data and School Quality Report data. To complete clustering, we replace all object data with numeric values, (eg Supportive Environment Rating: {'Not Meeting Target': 1, 'Approaching Target': 2, 'Meeting Target': 3, 'Exceeding Target': 4}), then added dummy variables, and normalized the dataset.

For the clustering part, we first try different numbers of clusters to fit all the data, and use the silhouette score to choose the number of clusters K. We used k=2 for clustering, then visualized all features by cluster0 and cluster1 and observed the difference in their numerical distribution. It can be seen that the following columns have obvious distribution differences:

Enrollment: The data in cluster0 is distributed from 0 to 800, and the data in cluster1 is concentrated from 0 to 2000.

Percent students with disability: The values of cluster0 are mostly 0.1 to 0.3, and the values of cluster1 are more evenly distributed from 0.1 to 0.5, and most of them are distributed from 0.25 to 0.4.

Percent in temp housing: The value of cluster0 is the highest at around 0.02, and then gradually declines. The temp housing percentage of cluster1 first rises and then falls, with a peak at 0.2.

Years of principal experience at this school: In cluster0, principal experience years range from 0 to 20, and in cluster1 the range is around 5 to 20.

Then we selected different subsets of the original dataset and divided them into: Basic data, education quality data, economic data, race data and without race data for clearer clustering and exploration.

Criteria for the classification is based on the correlation heatmap and following standards:

Subsets Data	Description	Columns
Basic data	Describing the most basic information about a school	Rigorous Instruction(%), Collaborative Teachers(%), Supportive Environment(%), Effective School Leadership(%), Rigorous Instruction Rating, Collaborative Teachers Rating, Supportive Environment Rating

Education Quality data	The dataset containing all the information describing the quality of teaching in schools	Years of principal experience at this school, Percent of teachers with 3 or more years of experience, Student Attendance Rate, Percent of Students Chronically Absent, Teacher Attendance Rate, Effective School Leadership Rating
Economic data	Describe the financial situation of the student's family (whether it is poor or not)	Strong Family-Community Ties(%), Trust(%), Percent Students with Disabilities, Percent Self-Contained, Economic Need Index, Percent in Temp Housing, Percent HRA Eligible, Strong Family-Community Ties Rating, Trust Rating
Race data	Data related to racial proportions of school students	Percent Asian, Percent Black, Percent Hispanic, Percent White
Without Race data	Remove all race information from the original dataset	All the columns except: Percent Asian, Percent Black, Percent Hispanic, Percent White

Cluster by Basic data (get 6 clusters):

Supportive Environment Rating: The values of cluster0 to cluster2 are small, basically distributed between 0 and 2.8, while the ratings of cluster3 to cluster5 are distributed around 3 to 4.

Cluster by Education Quality data (get 2 clusters):

Years of principal experience at this school: There are obvious differences between the data of the two clusters. For cluster 0, everyone has only 0 to 5 years of experience, while cluster1 has principal experience from 3 to 25 years.

Cluster by Economic data (get 6 clusters):

Trust(%): Only cluster0 and cluster2 are below 0.8.

Strong Family-Community Ties Rating: In this indicator, the ratings of cluster2 and cluster4 are all below 2.3, and the values of other clusters are mostly above 2.5.

Trust Rating: The trust ratings of cluster0 and cluster2 are all below 2.7, and above 2 in other clusters.

Cluster by race data (get 4 clusters):

Since the standard of clustering is race, we can see from the visualized column chart that only the four figures of Percent Asian, Percent Black, Percent Hispanic, and Percent White show significant differences in the numerical distribution trends. The distributions of the remaining columns are mostly similar.

Cluster by data without race (get 2 clusters):

Enrollment: The enrollment value in cluster0 is relatively large, ranging from 0 to 4000, and the value of cluster1 is around 0 to 700.

Percent students with disabilities: higher proportion of students with disabilities in cluster1.

Years of principal experience at this school: more principal experience years in cluster1.

In general, we can observe that under different clustering standards, there is basically no significant difference in the distribution of various statistical data of crime frequency. However, the data about students' family situation and school teaching quality show differences in distribution among different clusters, especially the indicator "Years of principal experience at this school", which can be focused on in follow-up research.

Methodology and Modeling

Random forest

Random forest is one of the most practical algorithms in bagging ensemble strategies. We chose this model because it has very high accuracy and is not easy to overfit. Furthermore, the random forest model can handle high dimensional data without dimensionality reduction.

For each category, we randomly divide the experimental data into two groups. 70% of the parcels were used to create the training area. The other 30% of the plots were used to create a test area to verify the accuracy of applying random forest to estimate crime performance. Then we fit the training dataset into GridSearchCV to optimize the RandomForestClassifier hyperparameter 'min_samples_split' and 'min_samples_leaf' between 2-10. Next we bring the best hyperparameter into RandomForestClassifier to predict the crime performance on the test dataset. Finally we calculate the accuracy of the model on each category, and the importance for every quality standards in each category.

Bayesian Network

Our dataset has a great number of variables and unclear relationships, thus Bayesian Network is a useful tool to provide a graphical representation of the structure of the probabilistic and causal relationship. It could both offer an understanding of interaction of variables and make predictions.

Our Bayesian network used the data after clustering in the cluster observation part. We had discretized the total crime number into two, crime, and non-crime in the data analysis part. Furthermore, we discretized 33 quality standards into 5 or fewer categories according to their distribution. And all the data are split into train and test datasets with a ratio of 8 to 2. For the PC algorithm, the data needs to be further turned into integers.

We use two algorithms, HillClimbSearch with ‘K2’ prior and BIC scoring function and PC algorithm, to learn the network structure, and use the learned structure to learn parameters, then make in-sample and out-of-sample predictions by passing values to the nodes in the network except for the binary crime node. In addition, since the PC algorithm provides the directed edges and undirected edges, we would try both to build the network.

For higher accuracy, the labels learned in the previous clustering part could be utilized to further divide the dataset into several sub-dataset (one sub-dataset for one cluster), and data from every cluster could be used to train a separate model. In our case, we found the direct parent, school type, of the binary crime node, has a big impact on it, thus dividing the dataset according to different school types is also experimented with. Table A1.0 shows our results.

Results

1 Results for Random Forest

For the performance evaluation, the random forest algorithm with the School Type has the highest mean accuracy of 0.86, followed by the one with financial situation of students, which had a mean accuracy of 0.773. The model with Education Quality and Race has the mean accuracy of 0.759 and 0.675. The model with the Basic category has the lowest mean accuracy, which is 0.623.

After that, we ranked the feature importance of the 33 quality standards according to the measurement results of the Random Forest Classifier. We learned that the supportive environment is the most important standard in the Basic category; student attendance rate is the most important standard in the Education Quality category; strong family-community ties is the most important standard in the Economy category; there is not much difference in importance between different races. Last but not least, in the School type category high school is the most important type that affects the variable of binary crime.

Fig.1 Heatmap



Fig.2 3 *Basic; Education Quality*

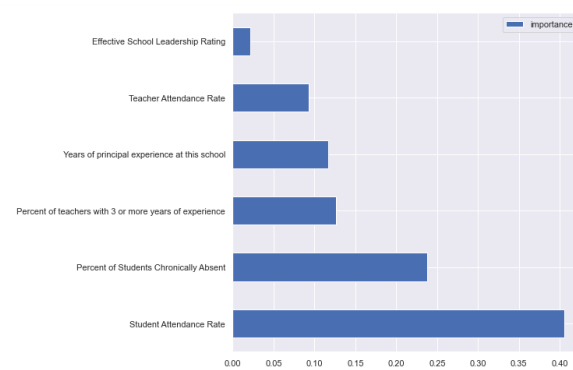
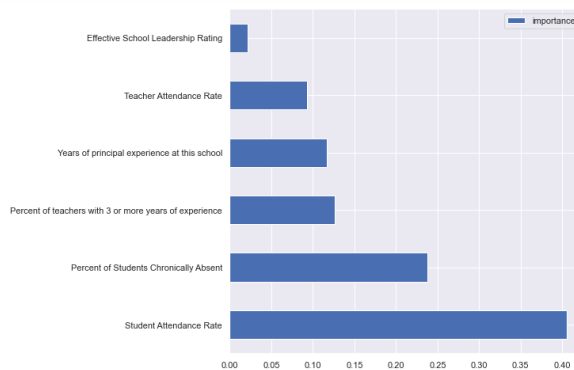


Fig.4 5 *Economy; Race*

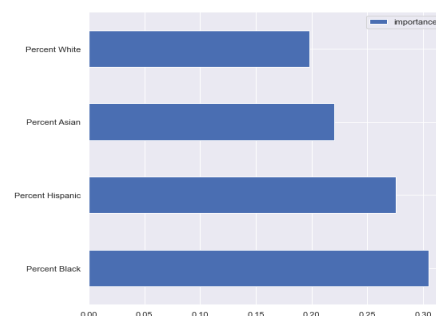
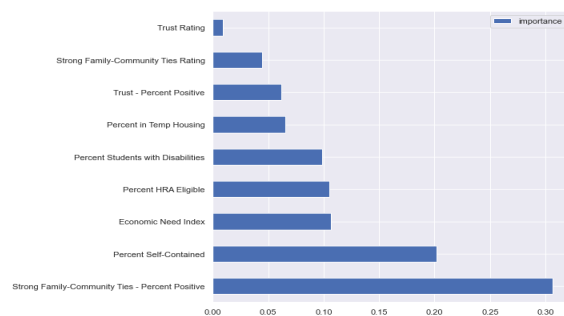
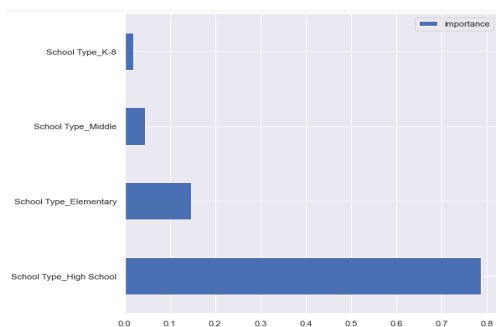


Fig.6 *School Type*



2 Results for Bayesian Network

2.1 Accuracy & Robustness

From Table A1.0, we could see the network learned by an undirected PC makes predictions with higher accuracy than HillClimbSearch. That is because HillClimbSearch may stop when finding local optima. When the local optima happen to be the global optima as in Education Quality Cluster_0, the two algorithms reach the same network structure and they have the same prediction accuracy. The performance of directed PC is better when testing the in-sample accuracy while worse in out-of-sample accuracy than HillClimbSearch, which may be due to the occasionality of the orientation of the remaining edges in the PC algorithm. However,

due to the limited amount of data in small clusters and a great number of explanatory variables, the PC algorithm sometimes fails while HillClimbSearch is still able to work well, which proves the robustness of HillClimbSearch algorithm in the case of limited data.

After dividing the whole dataset into sub-dataset based on clustering results, the accuracy changed. Ways of division based on different standards got different results. Divisions based on overall clusters and school types led to improved accuracy, while the performance of the model based on education quality clusters is even worse. Thus the first two ways of division are more effective for building high-accuracy models.

And the models for different clusters in one division standard also perform differently. For example, when dividing data based on school types, the prediction accuracy of the high school Bayesian network increases largely, the prediction accuracy of the Middle school network enhanced a little, while the prediction accuracy of the elementary school network decreased a little. That may be because the crime-related pattern in high school is much more explicit and predictable. While for K-8 schools, the binary crime node does not link to any other nodes thus the binary crime could not be inferred.

2.2 Interpretation

From Fig.D2.1-2.3 the Bayesian network for all the data, we could find school type and economic index are two key nodes that have high centrality. School type node link with nodes of attendance, enrollment to binary crime, strong family ties, self-contained, and a series of education and management related nodes. Economic need index node links races to the percentage of English learners and disabilities, and other economic situation nodes. And these two key nodes are linked by the chronicle absence node. Due to the apparent impact of school type on binary crime, we further analyze the networks under different school types. The networks show although there are detailed differences, the relationship chain and casual chain are similar in these networks as below:

- *Financial Situation of Students – Economic Need Index – Races – Percent English Learners*
- *Disability – Percent Asian*
- *Economic Need Index → Percent of Chronically Absent Student → Enrollment → Binary Crime*
- *Effective School Leadership → Collaborative Teachers → Positive Trust → Strong Family-Community Ties*
- *Effective School Leadership → Collaborative Teachers → Rigorous Instruction → Supportive Environment → Student Performance*

Thus Percent of Chronically Absent Students can directly affect binary crime, the economic situation of students has indirect influence for the binary crime. Effective school leadership could enhance the education quality and student performance, but has nothing to do with binary crime. Students of different races have different financial situations, English abilities, and disabilities, but races have no relationship with other issues.

However, according to Student Safety Act Reporting by NYCLU, it may be either the high Percent of Chronically Absent Students causes crimes to happen in school or the students in trouble are absent from school to avoid the police. And since the binary crime node is independent of the race node conditioned on the financial situations of students in all our networks, it could be indicated that compared to the race difference, the gap in wealth is more likely to be the cause of crimes in school. Though race difference is also the reason for unevenness of wealth, at least we could pay more attention to the gap in wealth so that it is not masked by race problems.

Discussion

Discussion for Random Forest

The performance of random forest classification will be affected by the noises and dependencies of the independent data in the model. In some of the categories, the features like supportive environment rating and supportive environment percent positive are not very independent. They thus will create noise and lower the accuracy of the model. Also, the distribution of our dependent variable total crime is unnormal, over 75% of the schools in most of the category datasets are without the crime records. This might be a reason for low accuracy as well.

Discussion for Bayesian Network

In Bayesian network analysis, categorized value is found much more valid to build the connected network than continuing value. When building networks with real-value data, the resulting network is fractal as Fig. D2.1, D2.2, and D2.3 show, and the prediction accuracy is low or not be able to infer since there are few or no nodes connected with the binary crime node. In addition, the inference power of Bayesian network analysis is constrained by the amount of data and whether the dependent variable node is linked with other nodes. In this case, the Bayesian network is invalid and the decision tree and the random forest are much more robust. Though the Bayesian network could learn the casual orientation structure from data, when combined with the reality we found it is not fully reliable, which may be because the assumption of PC algorithm, causal sufficiency is broken. Nevertheless, it is still a good way to understand the variable structure.

Discussion for Both Models

Finally, combining the results we get from Random Forest and Bayesian Network. We find that High School has the strongest influence on students' crime performance. We think this is because students in the stage of high school tend to stir up trouble. Meanwhile, K-8 has the lowest influence on students' crime performance, because K-8 has a relatively stable education environment.

In our case, Random Forest has higher prediction accuracy than Bayesian Network, and it is more robust to get the prediction results in the case of a limited number of data. However, Bayesian Network could compute probability distributions for unobserved variables but Random Forest could not. On the other hand, both models have powerful interpretability and are effective for dealing with large numbers of variables. Interestingly, the variable groups got from the close importance index of Random Forest and from the close network nodes are highly coincidental, which means a truly closer connection of variables within each group.

Conclusion

In this paper, we try to sort out the relationship between variables related to school quality and crime occurring in New York City public schools, and use variables from school quality to predict whether the crime happened there that year. First, we used clustering to explore the possible relationship. Then we innovatively used the labels resulting from clustering to learn separate Bayesian Network models for each group to improve the prediction accuracy. As a result, the prediction accuracy increased but differently in each cluster. The highest accuracy is over 0.81. We also divided directly the variables into several groups and built Random Forest models for each, and the result shows the highest accuracy is 0.86. The prediction power of the groups of variables:

School Type > Financial Situation of Students > Education Quality > Race > Basic Attributes of Schools

In conclusion, students' financial situation and school teaching quality show explicit distribution differences among different clusters. High schools have the most records of crime occurring in school, a high number of chronically absent students would directly increase the crime in school, and the bad economic situation of students has underlying implications for crimes. Effective school leadership could enhance the education quality and student performance, but has nothing to do with crime in school.

Contribution of each team member

Baohan Liu	Background, Data and Preprocessing, Cluster Code & Article
Fanshu Li	Abstract, Problem Definition, Cluster Code, Bayesian Code & Article, Conclusion

Jacob Jiang	Problem Definition, Data preprocessing, Random Forest Code & Article
Wangtianhan Pang	Background

Code written for the project and the datasets used can be found on:

<https://github.com/AsatsukiXiao/ML-Final-Group-Project.git>

Reference

[1] Grace Chen. (2022). Police Make Hundreds of Arrests at NYC Schools Last Year.

<https://www.publicschoolreview.com/blog/police-make-hundreds-of-arrests-at-nyc-schools-last-year>

[2] NYCLU. (2019). Student Safety Act Reporting.

https://www.nyclu.org/sites/default/files/ssa_sy_17-18_factsheet_nyclu.pdf

[3] NYCLU. (2018). New Data: Police Disproportionately Target Black And Latino Students In The NYC Schools

<https://www.nyclu.org/en/press-releases/new-data-police-disproportionately-target-black-and-latino-students-nyc-schools>

[4] Lankford, Hamilton, Susanna Loeb, and JamesWyckoff, “Teacher Sorting and the Plight of UrbanSchools: A DescriptiveAnalysis,”Educational Evaluation and Policy Analysis, 24 (2002), 37–62.

[5] Cook, Philip J., Denise C. Gottfredson, and Chongmin Na, ”School Crime Control and Prevention,” Crime and Justice, 39 (2010), 313–440.

[6] Deming DJ. Better Schools, Less Crime?. Quarterly Journal of Economics. 2011;126 (4):2063-2115

Appendix A

A1.0 Results of Bayesian Network

Data Division Standard	Cluster Type	Graph Index In Appendix	Network Learning Method	In-sample Accuracy	Out-of-sample Accuracy
All the Data	----	D 2.1	HillClimbSearch	0.67	0.68
		D 2.2	PC(directed)	0.68	0.64
		D 2.3	PC(undirected)	0.75	0.69
School Types	Cluster_0 K-8	A 1.1	HillClimbSearch	—	—
	Cluster_1 Elementary School	A 2.1	HillClimbSearch	0.62	0.66
		A 2.2	PC(undirected)	—	—
		A 2.3	PC(directed)	—	—
	Cluster_2 Middle School	A 3.1	HillClimbSearch	0.78	Failed
	Cluster_3 High School	A 4.1	HillClimbSearch	0.81	0.87
Education Quality Clusters	Cluster_0	B 1.1	HillClimbSearch	0.66	0.73
		B 1.2	PC(undirected)	0.66	0.73
		B 1.3	PC(directed)	0.66	0.73
	Cluster_1	B 2.1	HillClimbSearch	0.66	0.68
		B 2.2	PC(undirected)	—	—
		B 2.3	PC(directed)	—	—
Overall Clusters	Cluster_0	C 1.1	HillClimbSearch	0.65	0.72
		C 1.2	PC(undirected)	0.76	0.65
		C 1.3	PC(directed) Algorithm	0.72	0.71
	Cluster_1	C 2.1	HillClimbSearch	0.79	0.69

* “----” means the “Total N” binary crime variable are not a node in this network, thus could not be predicted.

* Blue numbers mean they are the same as above.

* The results of failed PC algorithm due to insufficient data are not listed

Graphs For Bayesian Networks of School Types

Fig. A1.1



Fig. A2.1

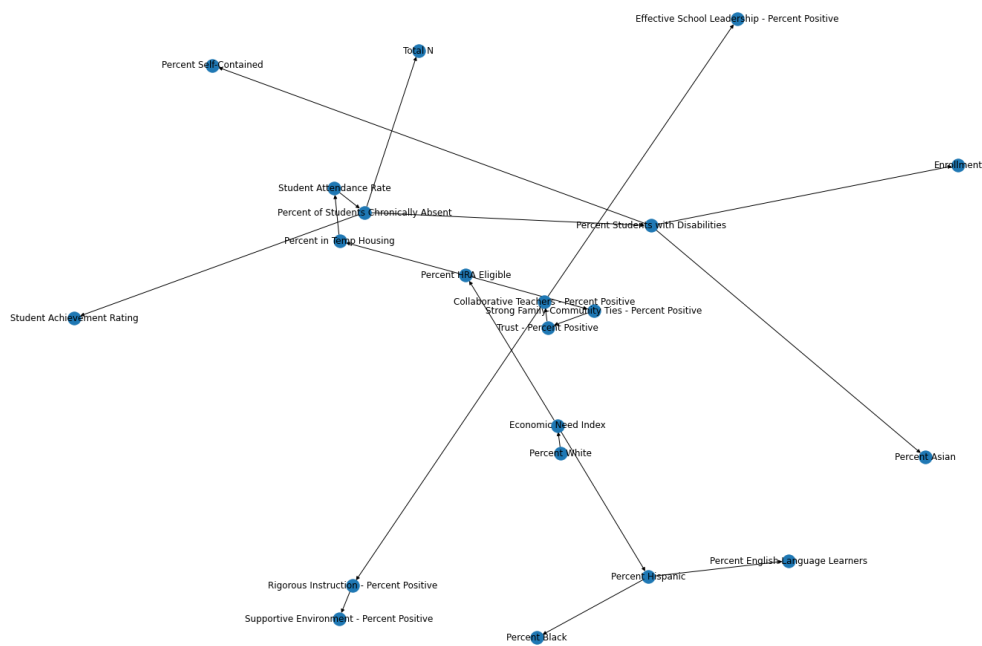


Fig. A2.2

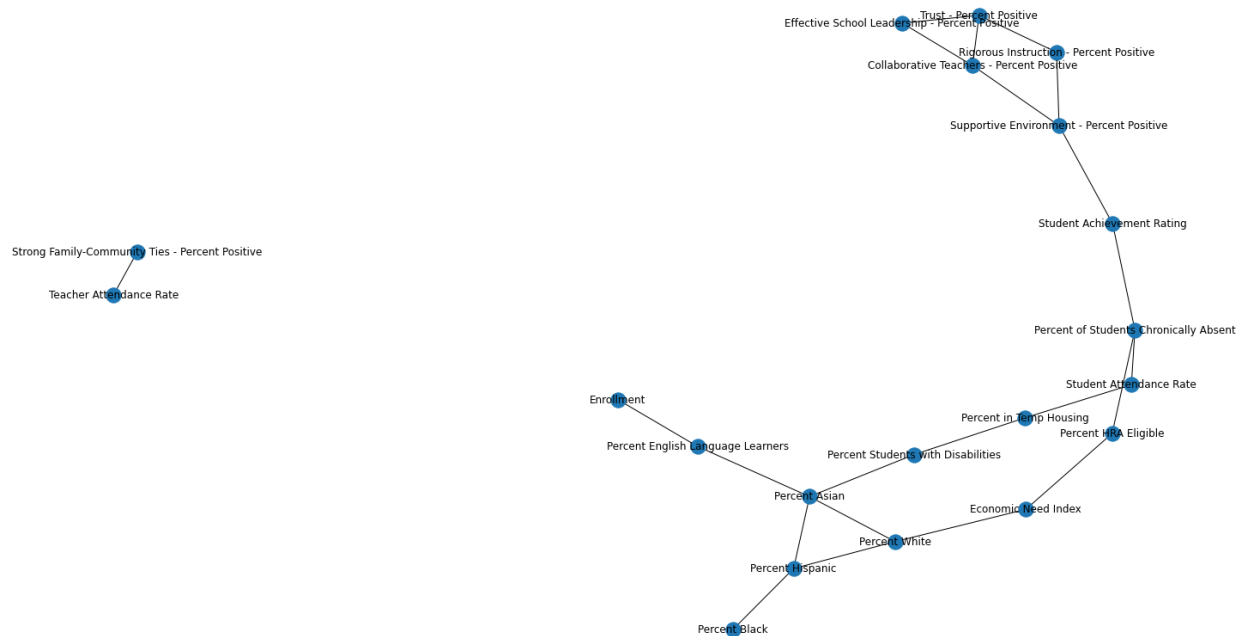


Fig. A2.3

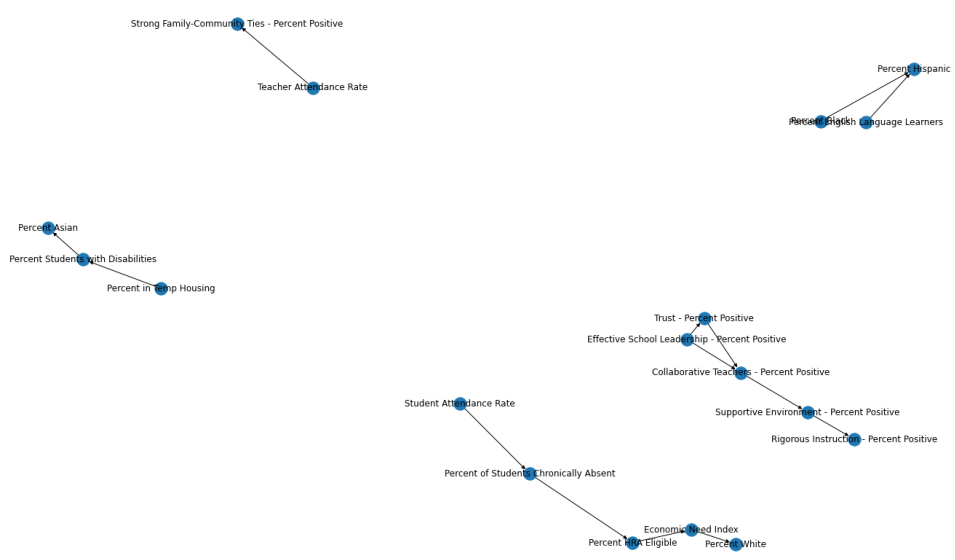


Fig. A3.1

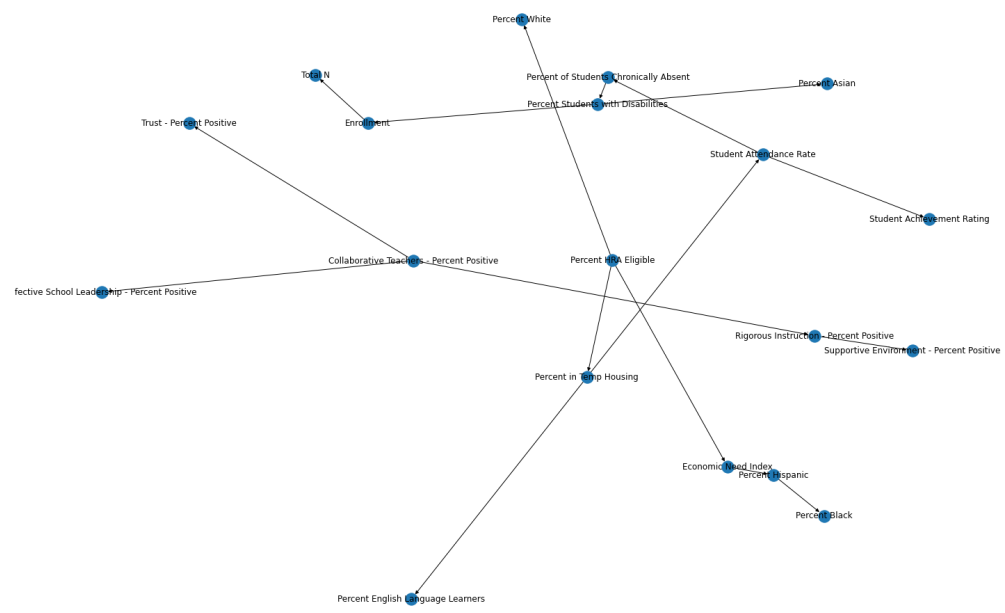
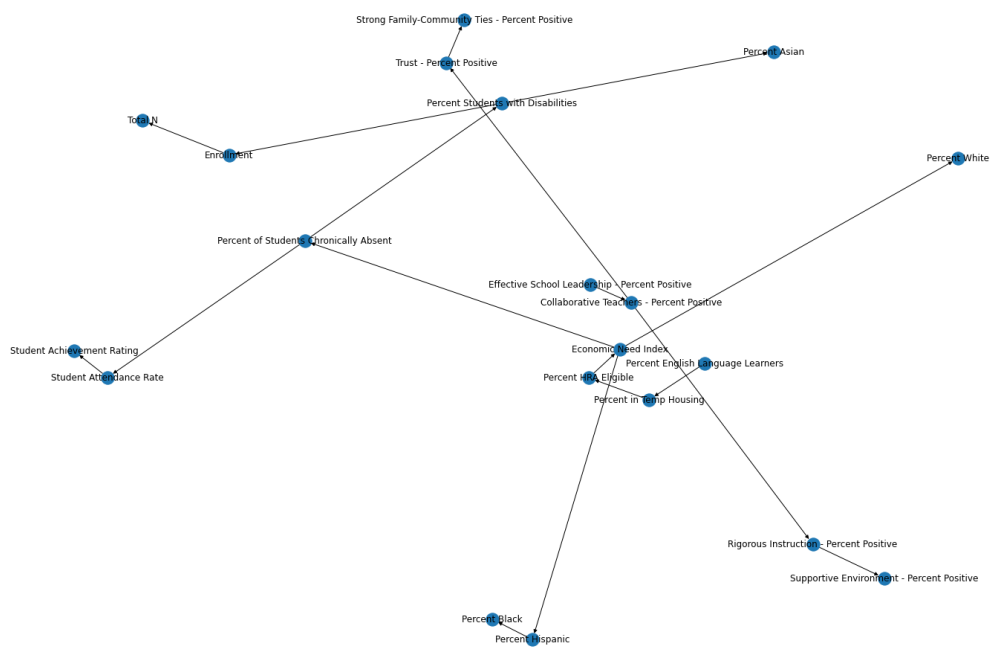


Fig. A4.1



Appendix B

Graphs For Bayesian Networks of Education Quality

Figure B 1.1

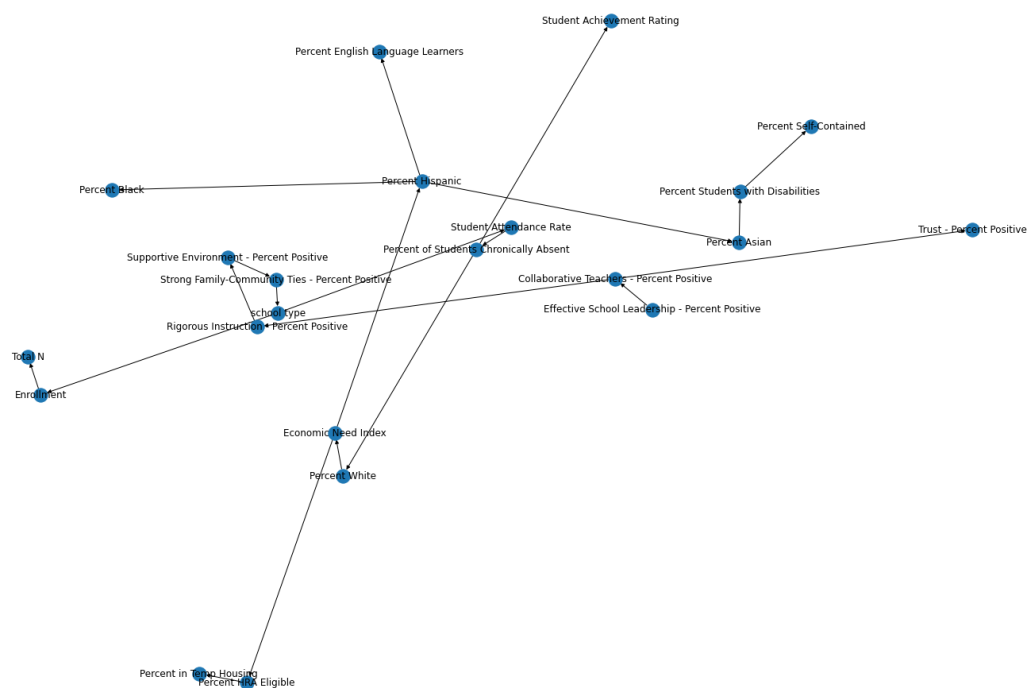


Figure B 1.2

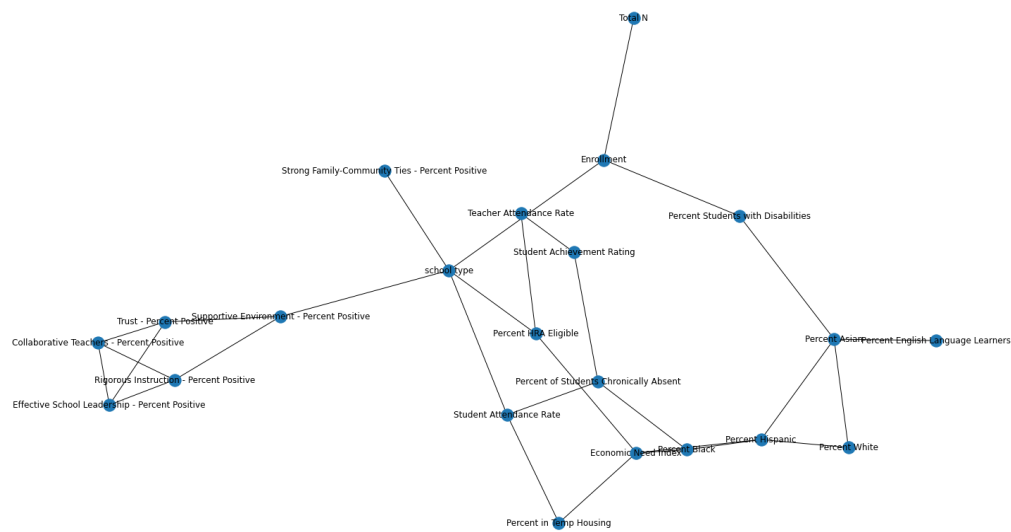


Figure B 1.3

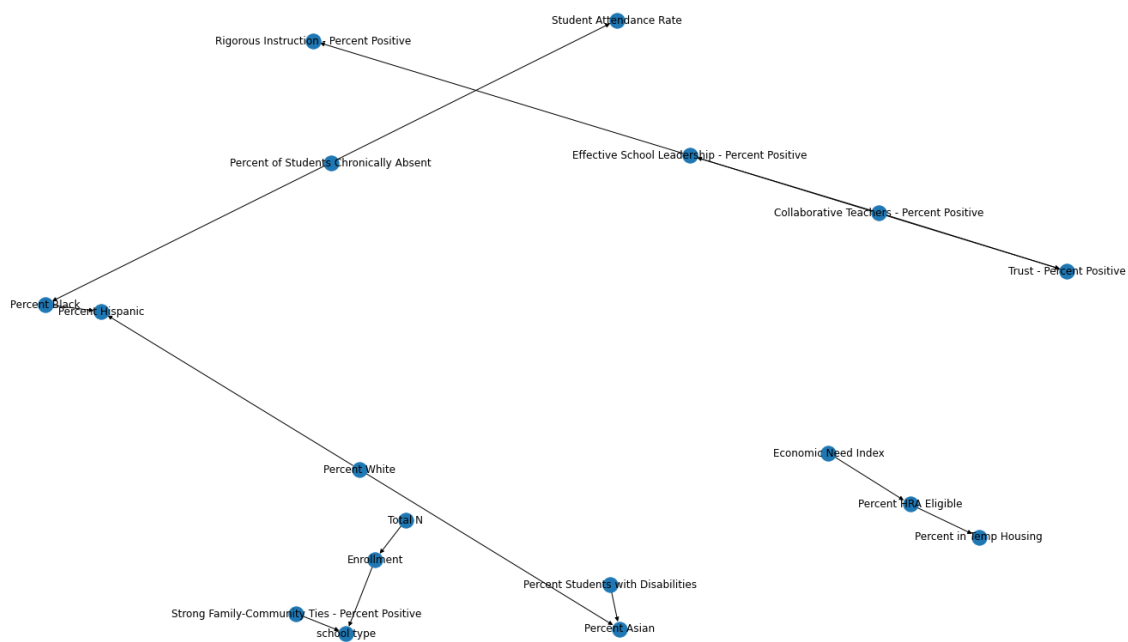


Figure B 2.1

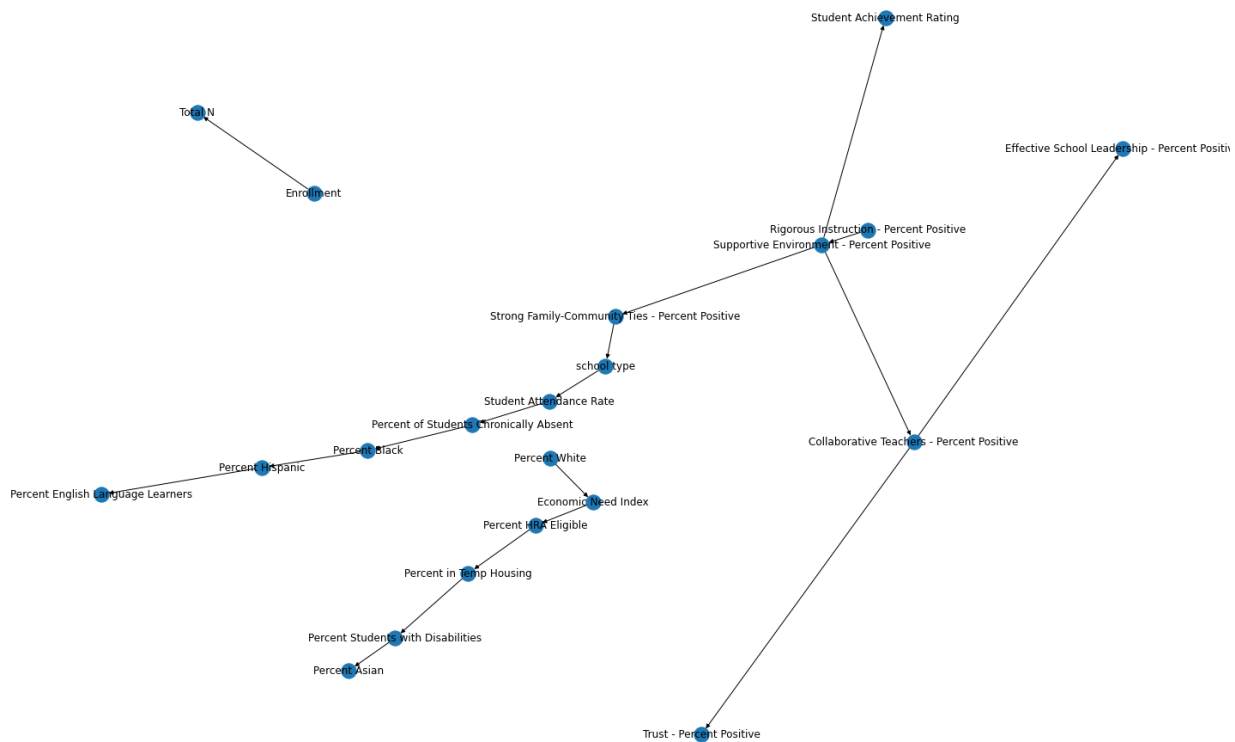


Figure B 2.2

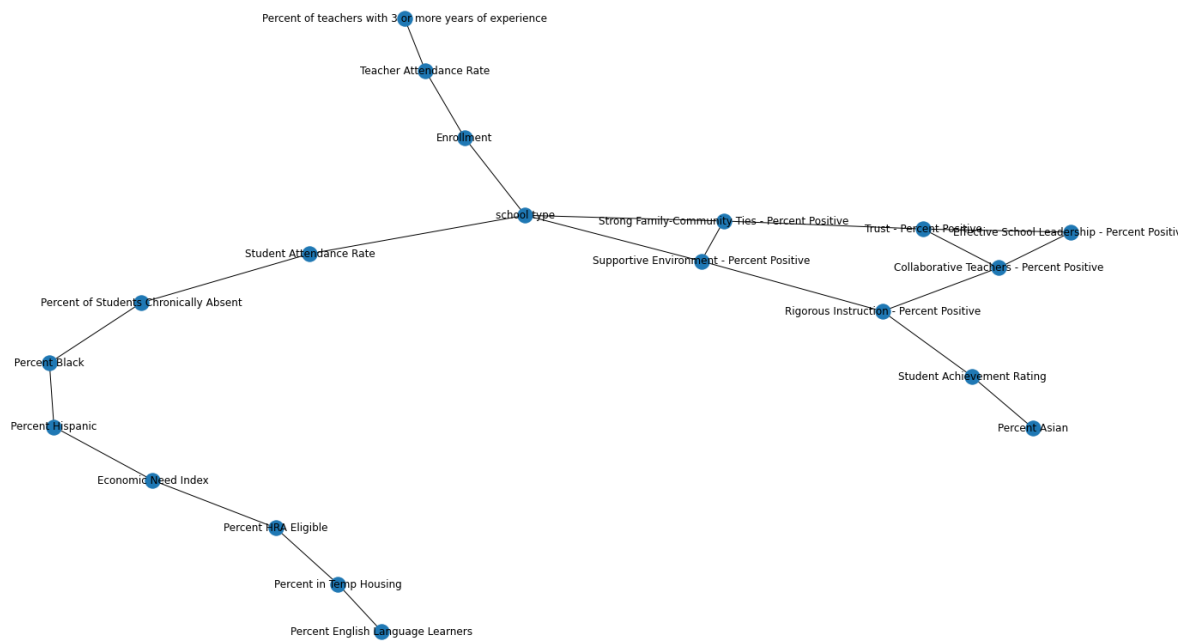
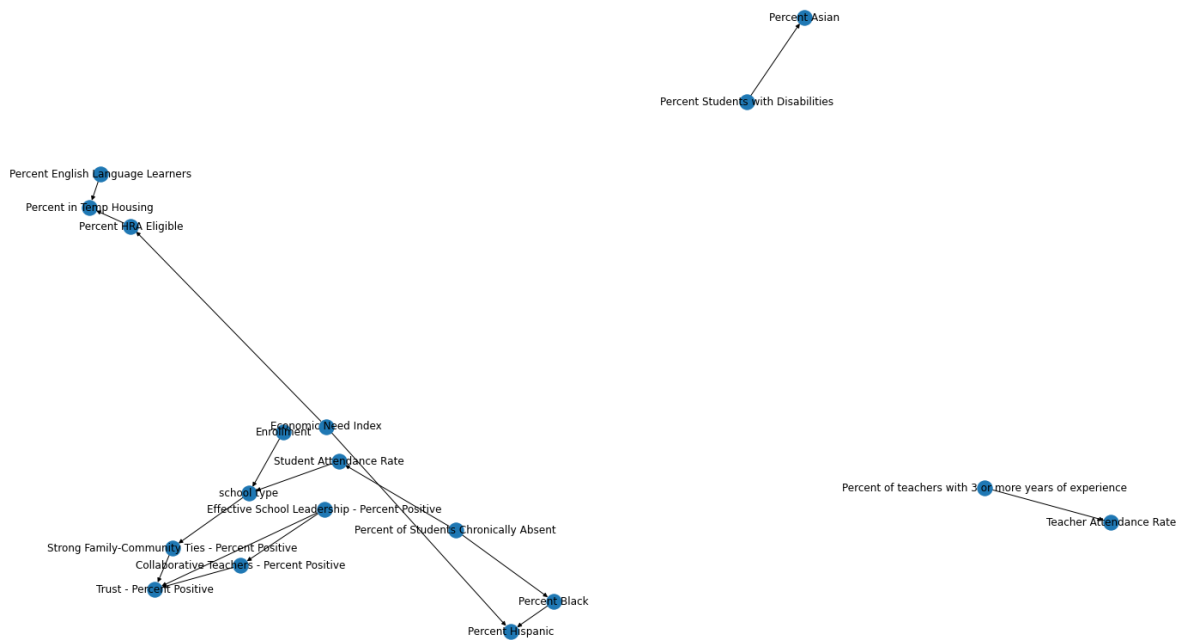


Figure B 2.3



Appendix C

Graphs For Bayesian Networks After Overall Clusters

Figure C 1.1

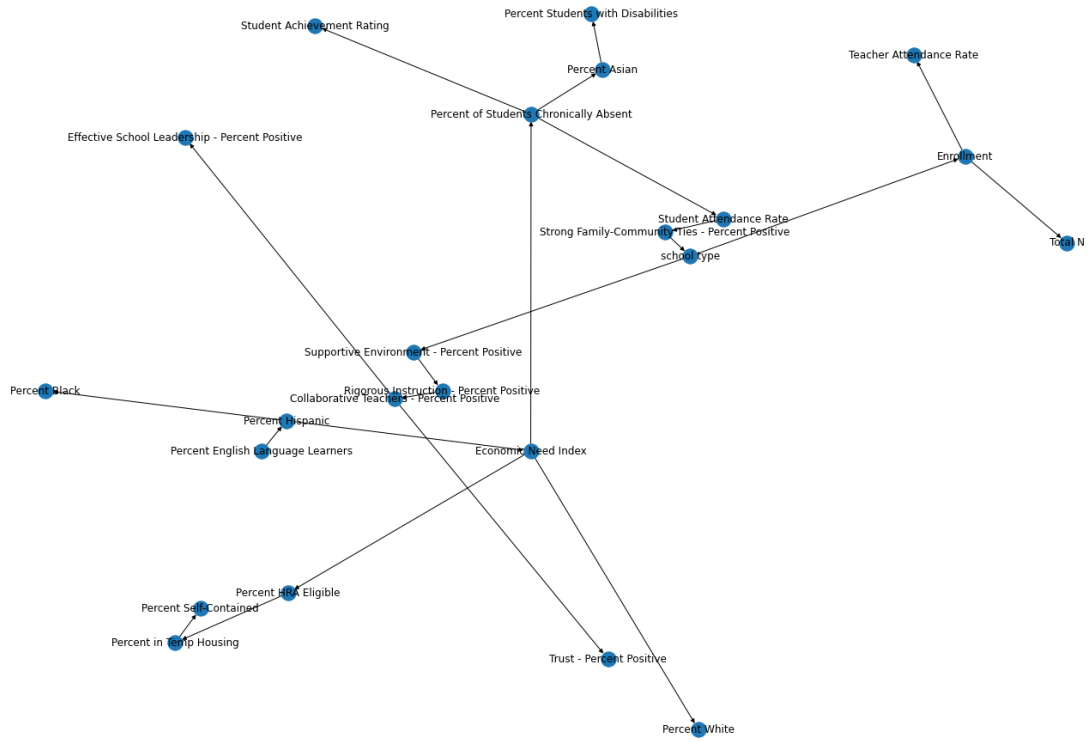


Figure C 1.2

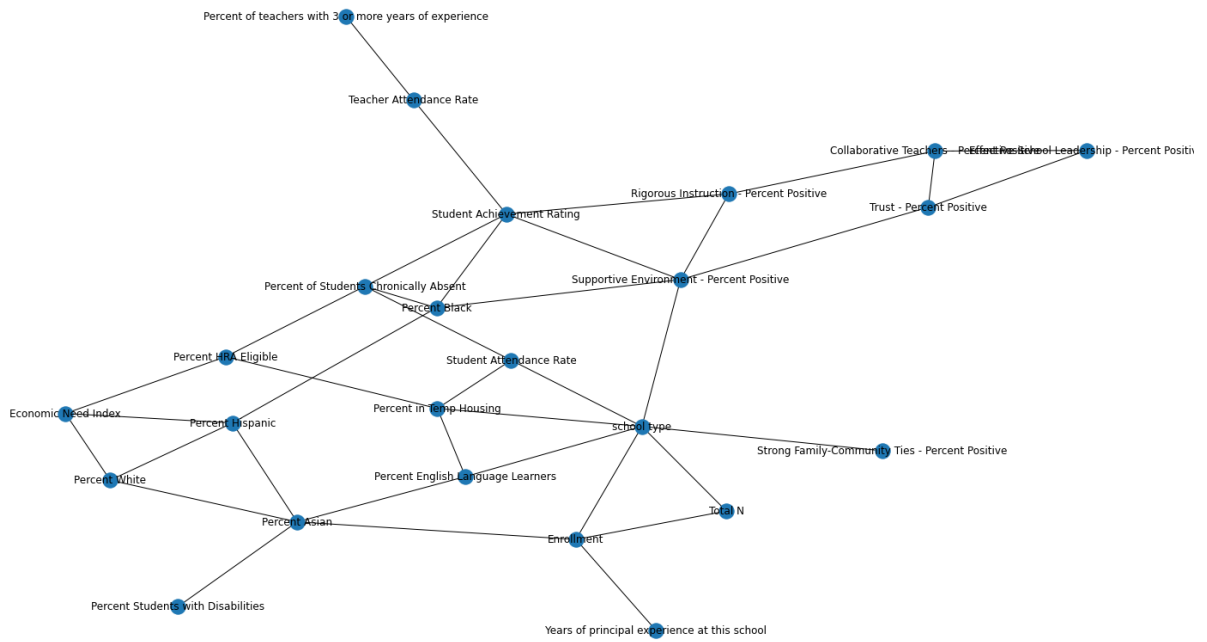


Figure C 1.3

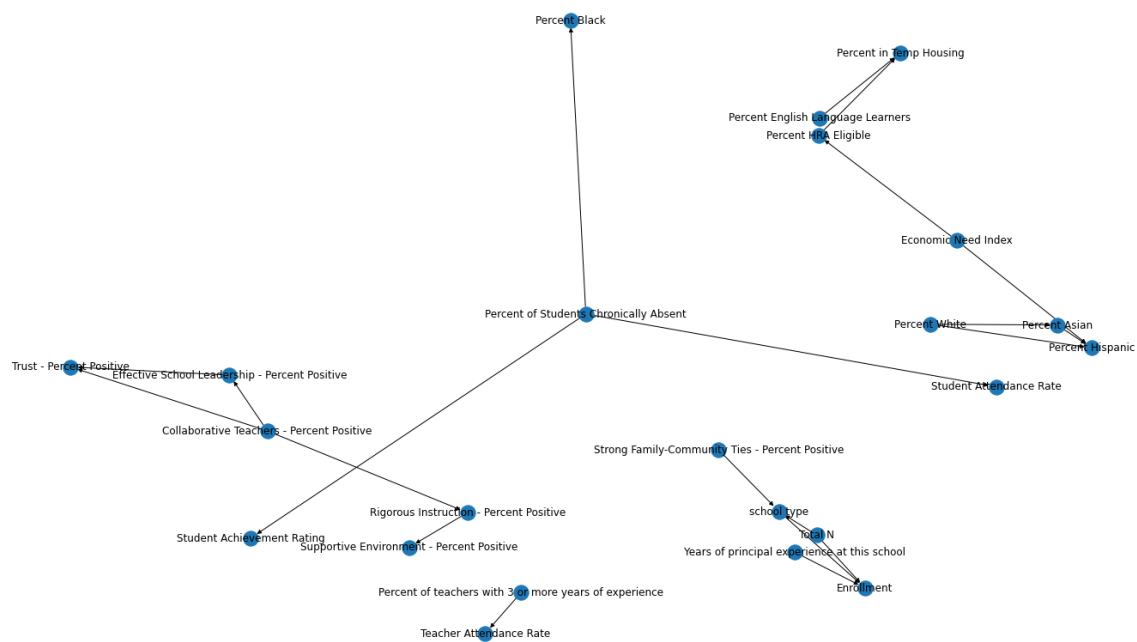
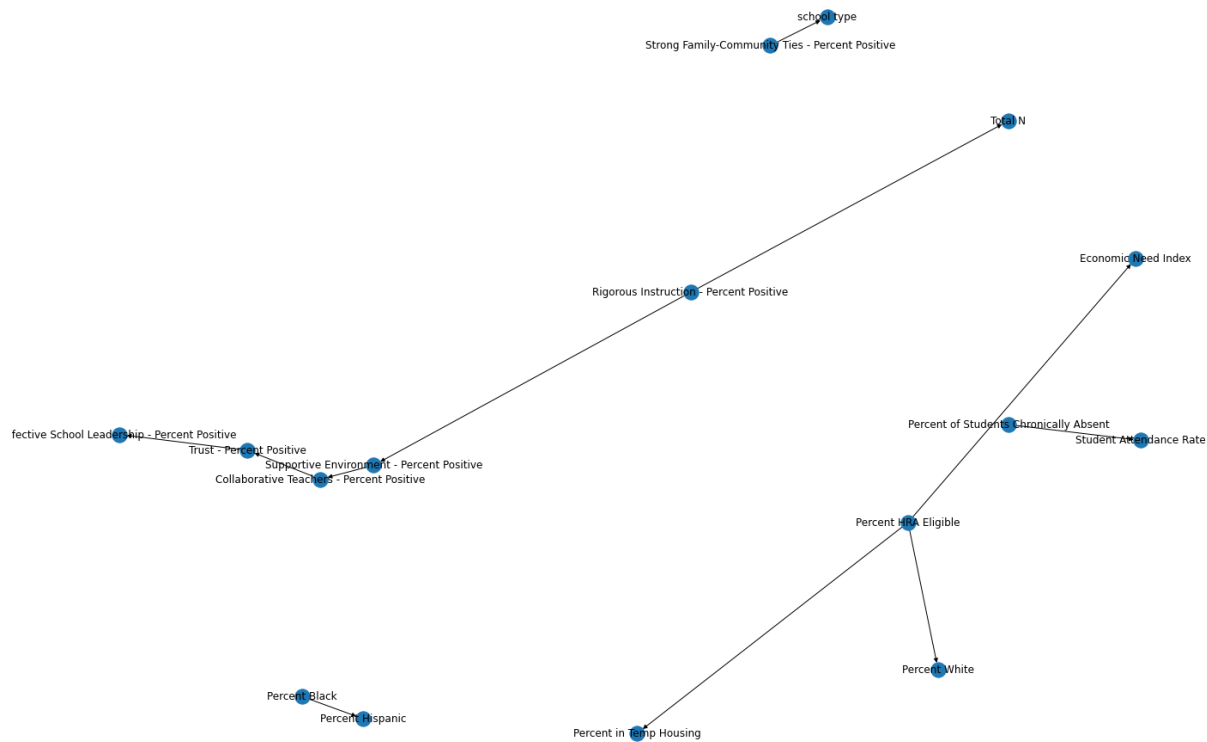


Figure C 2.1



Appendix D

Fig. D 1.1

Real value - HillClimbSearch - All the data

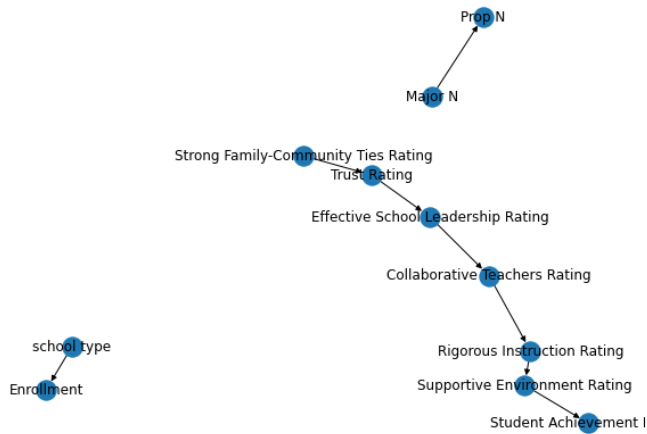


Fig. D 1.2

Real value - PC(undirected) Algorithm - All the data

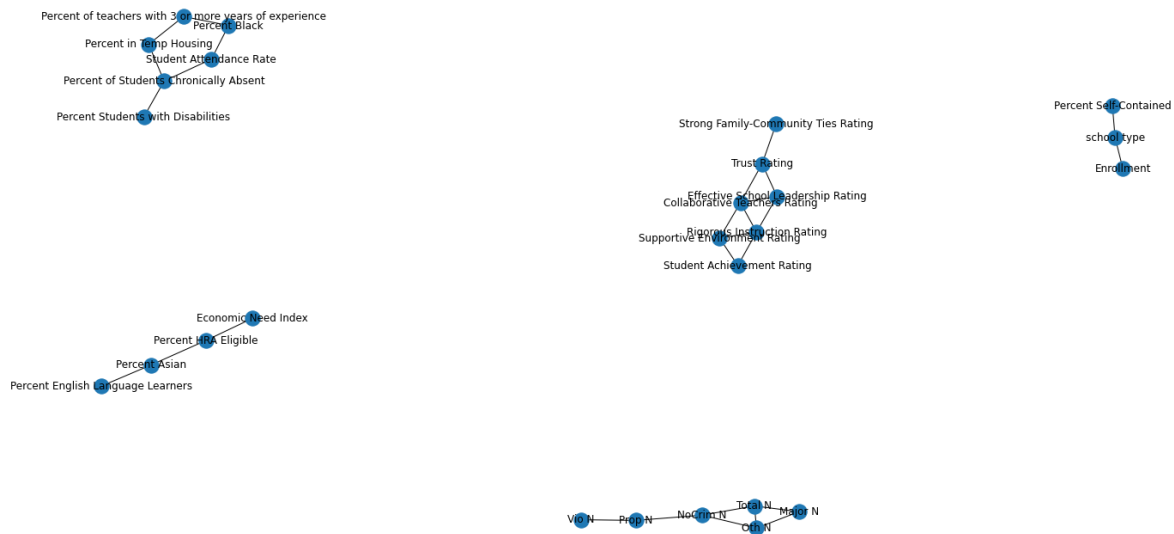


Fig. D 1.3

Real value - PC(directed) Algorithm - All the data

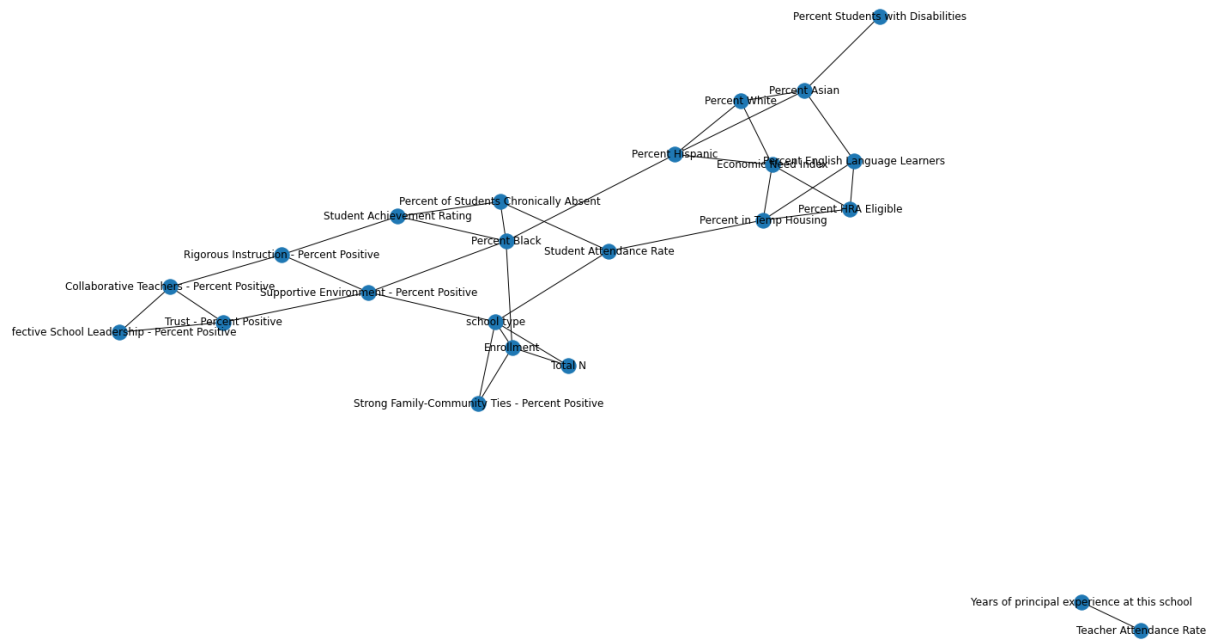
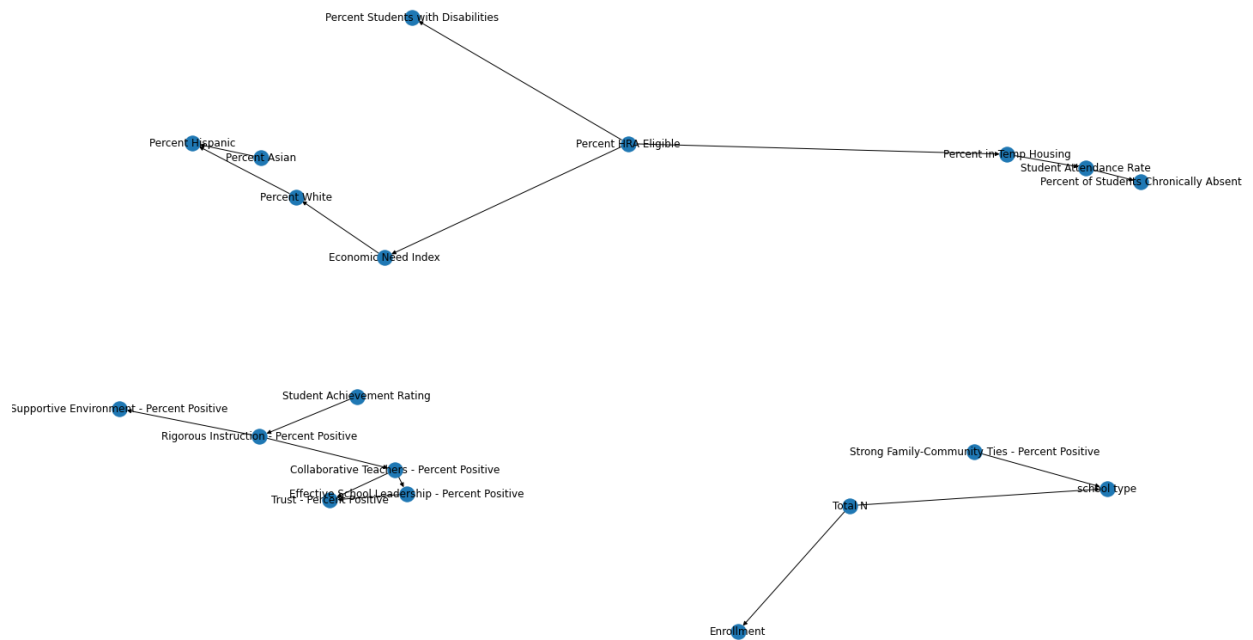


Fig. D 2.3



Appendix E

Graphs For Bayesian Networks After Student Financial Situation Clusters

Figure E 1.1

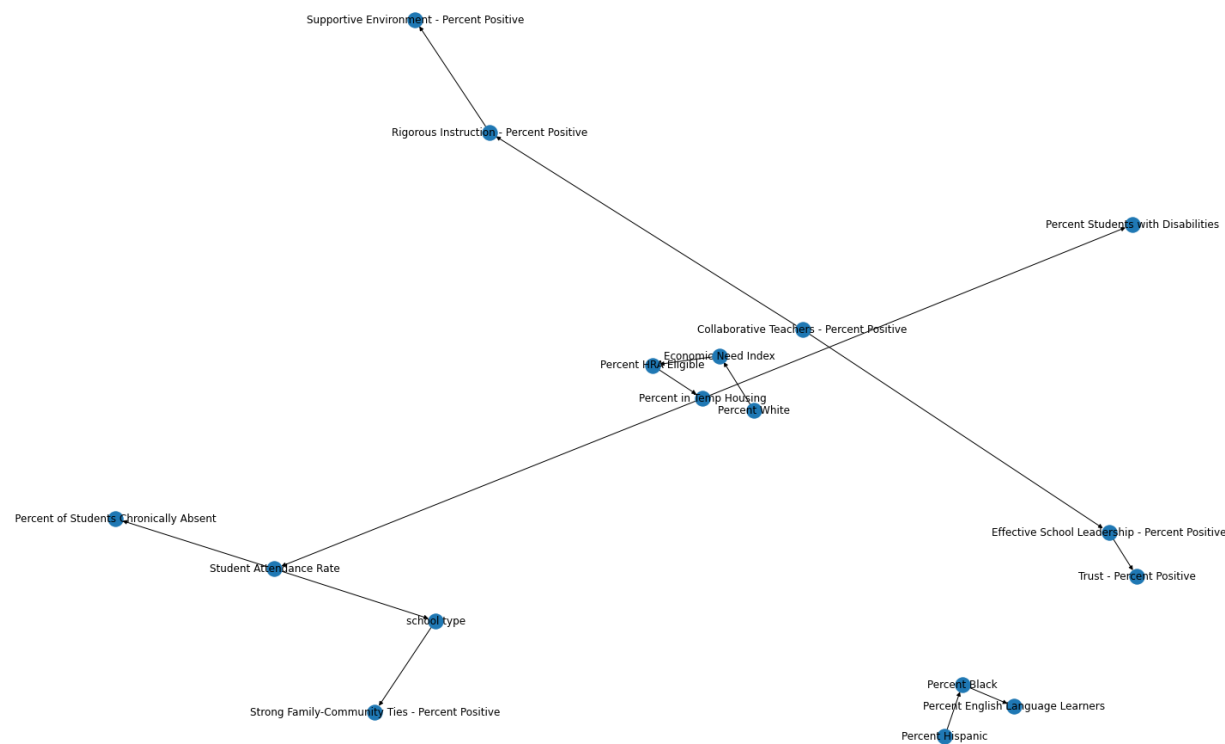


Figure E 2.1

Figure E 2.3

