Coursera Applied Data Science Capstone Project

# The Battle for Schools

Ruolin Fan

## 1. Introduction

- Hong Kong is a city with extremely high population density and diverse culture. Different districts may possess very different features and environment. High gap of wealth is also a characteristic of this society, which can potentially bring along many social problems.

- Education has always been a hot topic and attracts a lot of attentions in HK society. Both students and parents face high competitive pressure. There are different types of primary schools with different funding agencies. While the educational resources and cultural influences may differ in these different types of schools, the neighborhood around these schools may also have certain patterns that might be influential to the students.

- In an attempt to focus on education in the context of Hong Kong districts, the following questions/aims are raised for this project:

    A. Are there any general features for venues nearby primary schools? Can these features help group schools into clusters?
    B. If the schools can be clustered, does it correlate to its category (government, private, international, etc.?), or the district where it is located?
    C. Is educational resource correlated with other features of the district, e.g. household income, rent, etc.?
    D. From the above observations, generate implications for students, parents or educators

## 2. Data

- List for primary schools with classified types and location

    A dataframe of 567 primary schools was obtained containing the latitude and longitude information as well as features including the dataset it belongs to (Aided, private, government, international, direct subsidy scheme and English Schools Foundation), the gender of students (boys, girls and co-ed), the religion (Nonreligious, Protestantism / Christianity, Catholicism, Buddhism, Taoism, Confucianism, Confucianism-Buddhism&Taoism, Islam and others) and the session information (whole-day and A.M.).

    Source: https://geodata.gov.hk/gs/

- Ranking of 100 primary schools based on multiple factors.

  Source:
  https://www.professionaltutor.hk/article/%E5%B0%8F%E5%AD%B8%E6%8E%92%E5%90%8D100%E5%BC%B7-top100-primary-schools-%E9%A6%99%E6%B8%AF%E5%B0%88%E6%A5%AD%E5%B0%8E%E5%B8%AB%E6%9C%83-professionaltutorhk-%E4%B8%8A%E9%96%80%E8%A3%9C%E7%BF%92-%E5%90%8D%E6%A0%A1%E5%B7%A1%E7%A6%AE

  The ranking for primary schools found on this website is based on a set of scores for multiple factors, including 'Academic', 'PE', 'Music', 'Teacher' and 'Housing'. Speculating from the original data presentation, these were weighted by 40%, 15%, 15%, 15% and 15%, respectively, to generate the value of "Score_this_year"; while the overall rating considering historical records is represented by "Score_overall". There is an additional column of "ranking" for the 100 schools evaluated by this website.

- Tuition of 54 primary schools.

  Source: https://www.myschool.hk/primary-school/Primary-School-Fee.php

- Area of 18 districts

  Source:
  https://opendata.esrichina.hk/datasets/eea8ff2f12b145f7b33c4eef4f045513_0/data?orderBy=ENAME

- Statistics for each district including median monthly domestic household rent, median monthly domestic household income, median mortgage payment and loan repayment, median rent to income ratio, etc.

  Source: https://data.gov.hk/en/geospatial-data

3. **Methodology**

**3.1 K-means clustering**

K-means clustering is a method used to classify data into multiple clusters. Here, the primary schools are clustered using this method with the type of their most common nearby venues as features.

**3.2 Dython package for association analysis**

In order to explore the correlations between different features of primary schools or the 18 districts, the "Dython" package developed by shakedzy was used

http://shakedzy.xyz/dython/. This package uses different analysis of correlation for different types of data:

• Pearson's R for continuous-continuous cases

• Correlation Ratio for categorical-continuous cases

• Cramer's V or Theil's U for categorical-categorical cases

The "association" function is used, which returns a dataframe containing correlation statistics corresponding to the respective datatype and generates a heatmap for visualization.

Here, Theil's U was used for categorical-categorical cases since its unsymmetrical feature reveals additional information for one-directional correlations.

## 4. Results

### 4.1 Find nearby venues of primary schools and clustering

The nearby venues with parameter radius = 500, limit = 100 was extracted for each primary school using Foursquare based on the location information. Only 559 out of 567 schools got the nearby venue information, which were retained in the later analysis (others were excluded). These venues consist of 300 unique categories, and the most common nearby venue categories were restaurants or food shops (Table 1). This might simply reflect a general enrichment for such venues in the neighborhood where these schools are located.

| School | |
|---|---|
| **Venue Category** | |
| Chinese Restaurant | 898 |
| Fast Food Restaurant | 745 |
| Coffee Shop | 587 |
| Café | 490 |
| Shopping Mall | 417 |
| Japanese Restaurant | 373 |
| Hong Kong Restaurant | 355 |
| Dessert Shop | 354 |
| Noodle House | 350 |
| Cha Chaan Teng | 336 |

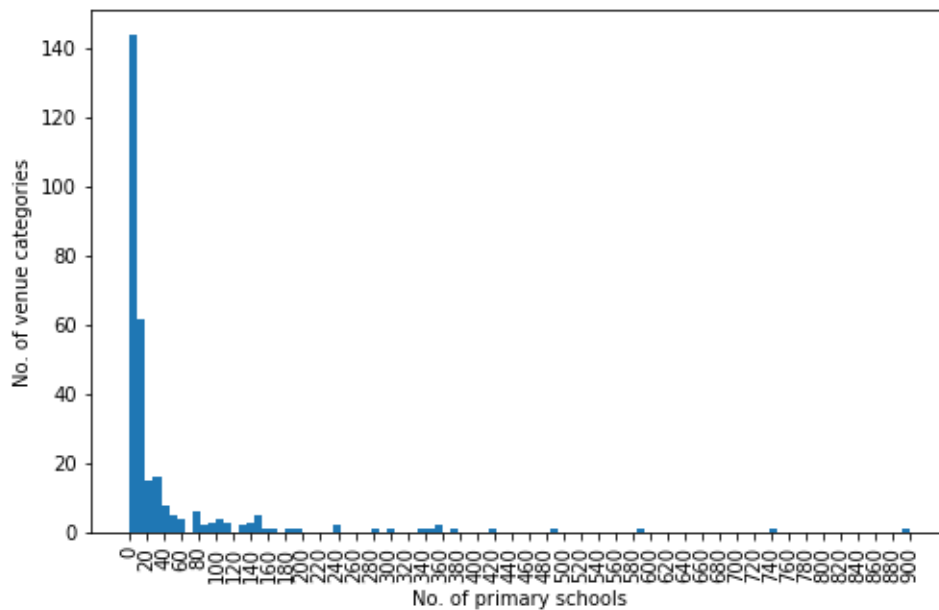**Table 1. The top 10 venue categories that are most commonly nearby primary schools.**



**Figure 1. Histogram for the venue categories identified to be nearby a certain number of primary schools.**

The top 10 most common venue categories for each school were summarized and used as features for clustering. The number of clusters was arbitrarily determined as 7. A cluster label from 0 to 6 was assigned to each school and the most common 10 venues among the "top 10s" of each cluster were summarized

in Table 2. The schools are also mapped with color coding for different clusters (Figure 2.)

| Cluster Labels | 1st Most Common Venue | 2nd Most Common Venue | 3rd Most Common Venue | 4th Most Common Venue | 5th Most Common Venue | 6th Most Common Venue | 7th Most Common Venue | 8th Most Common Venue | 9th Most Common Venue | 10th Most Common Venue |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Chinese Restaurant | Fast Food Restaurant | Hong Kong Restaurant | Shopping Mall | Asian Restaurant | Food Court | French Restaurant | Food Truck | Fried Chicken Joint | Food & Drink Shop |
| 1 | Fast Food Restaurant | Chinese Restaurant | Shopping Mall | Hong Kong Restaurant | Park | Cha Chaan Teng | Bus Station | Asian Restaurant | Convenience Store | Food Court |
| 2 | Food Court | Food & Drink Shop | Park | Food Truck | Flower Shop | French Restaurant | Chinese Restaurant | Fried Chicken Joint | Fujian Restaurant | Garden |
| 3 | Shopping Mall | Light Rail Station | Fast Food Restaurant | Fried Chicken Joint | Food Court | Flower Shop | Food Truck | Food & Drink Shop | Cha Chaan Teng | Flea Market |
| 4 | Fast Food Restaurant | Shopping Mall | Hong Kong Restaurant | Cha Chaan Teng | Bus Station | Light Rail Station | Chinese Restaurant | Market | Park | Food Court |
| 5 | Coffee Shop | Chinese Restaurant | Café | Food & Drink Shop | Japanese Restaurant | Food Court | Fast Food Restaurant | Food Truck | Flower Shop | Shopping Mall |
| 6 | Chinese Restaurant | Coffee Shop | Café | Fast Food Restaurant | Dessert Shop | Cantonese Restaurant | Noodle House | Shopping Mall | Japanese Restaurant | Cha Chaan Teng |

**Table 2. The 10 most common venue categories in each cluster of schools.**

An instant observation is that a majority of the venues are restaurants or food shops, same as the previous analysis of all schools. This indicates that the clusters might be essentially similar to each other. This makes it hard to deduce much relevant information simply from this table.

```
6    242
1    149
4     67
5     56
2     28
0     12
3      5
Name: Cluster Labels, dtype: int64
```

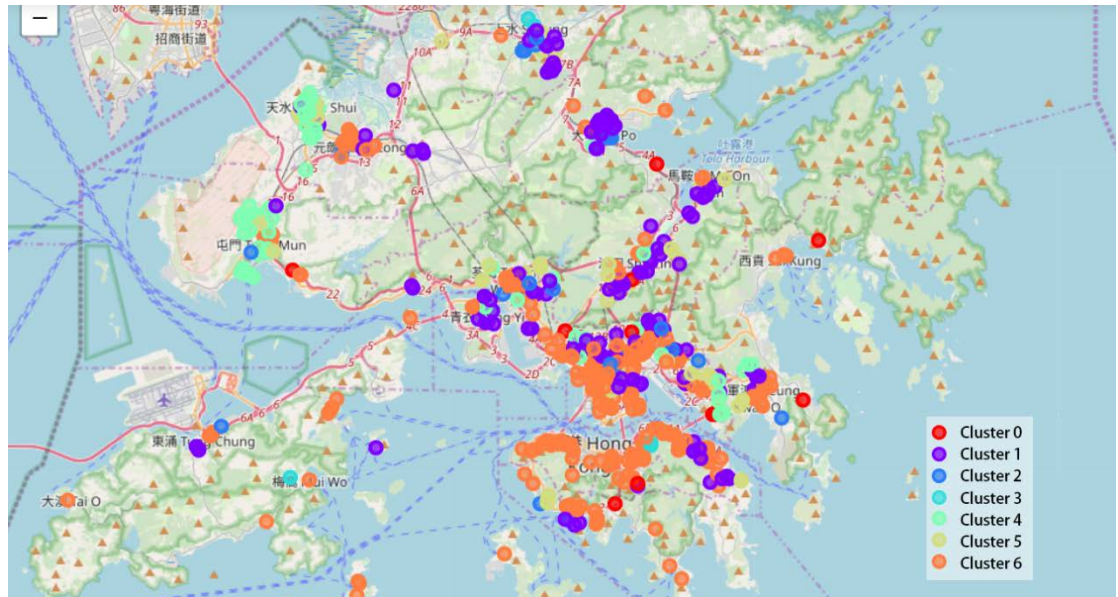**Figure 2. Frequency table for primary schools belonging to each cluster.**

**Figure 3. Distribution of primary schools in each cluster of common nearby venues.**

Value counts of individual cluster labels reveals that cluster 6 and cluster 1 were the largest 2 clusters, containing almost 70% of the schools. Interestingly, from the map it could be observed that the different clusters seem to be also segregated geologically. Cluster 6 schools predominantly occupies Central and Western district as well as Kowloon district; while cluster 1 is expands over New Territory. Considering that cluster 1 contains bus station and park in the most common venues as compared to cluster 6 which basically only contains food stores, this is probably because the Central and Western District and Kowloon areas are more commercial and populated areas with relatively more restaurants as compared to transportation points or open areas (park). Cluster 4 is predominantly localized in Tuen Mun and Yuen Long, which is also reasonable since the light rail station is one of the most common nearby venues. This is also distinct from cluster 3, which is much less in number, probably because cluster 3 has even more light rail stations (2nd most common) than restaurants, which could be uncommon.

## 4.2 Other features of primary schools

How are schools with different features localized geologically? Is there a certain pattern in their localization? In order to explore this question, schools with different categories were mapped with distinct colors. Features explored include dataset, student gender and religion, while the session feature was excluded since all but one school is in the type whole-day。

**4.2.1 Dataset.**

```
Aided Primary Schools                      411
Private Primary Schools                     64
Government Primary Schools                  33
International Schools (Primary)             24
Direct Subsidy Scheme Primary Schools       19
English Schools Foundation (Primary)         8
Name: Dataset, dtype: int64
```

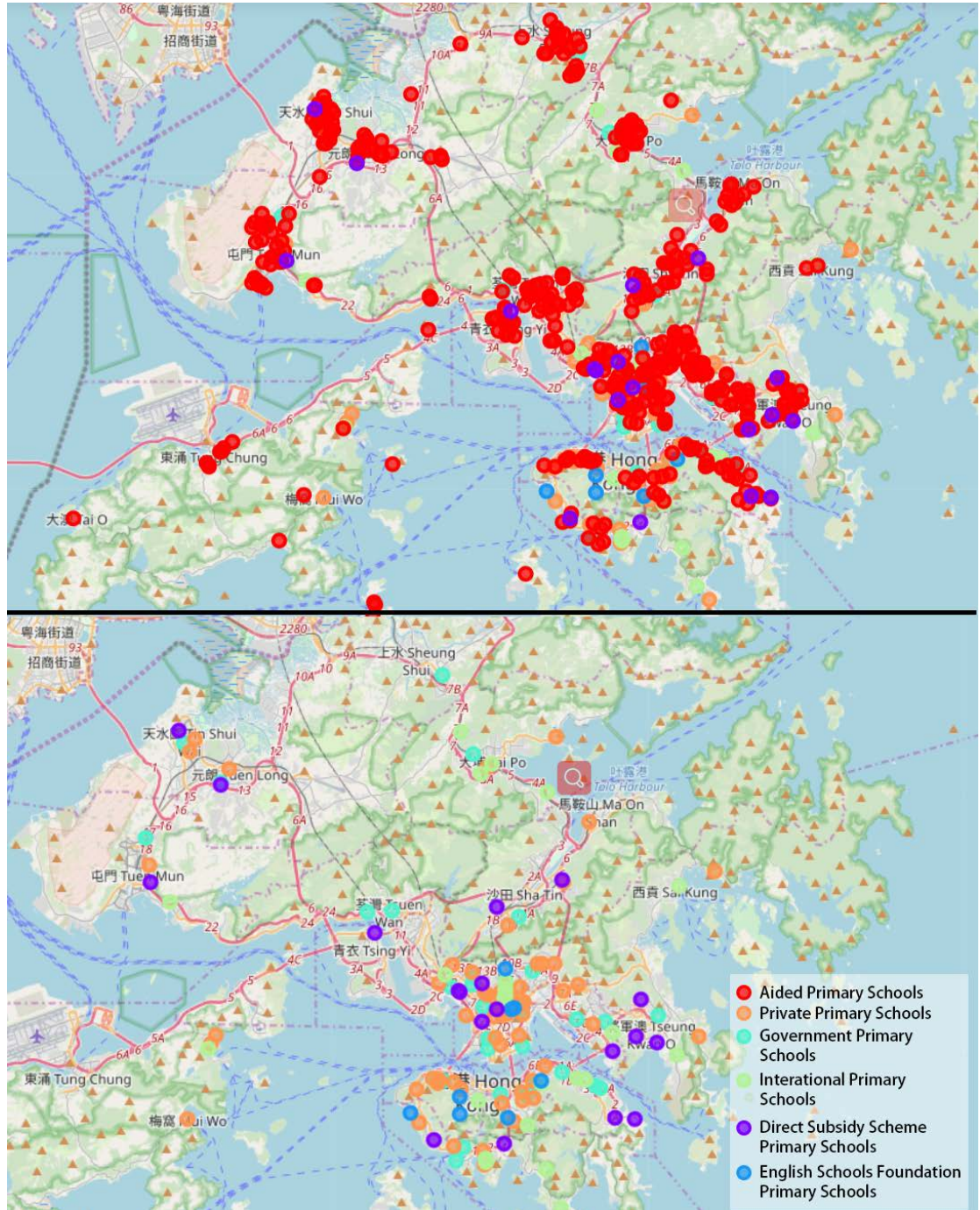**Figure 4. Frequency table for primary schools belonging to each category of school.**

**Figure 5. Distribution of primary schools belonging to each dataset (different category of schools). Upper: map with mark of all 6 categories; bottom: map with mark of all but the most common category (Aided Primary Schools).**

From the maps we could observe an overall dispersed and even distribution of each category of schools. The only exception is English Schools Foundation Primary Schools, which are restrained in Hong Kong Island and West Kowloon area. From a regional perspective, the primary schools around Sheung Shui and Tai Po area were relatively rare and limited to either Government or International Schools.

## 4.2.2 Gender

```
CO-ED      531
GIRLS       18
BOYS        10
Name: Students Gender, dtype: int64
```

**Figure 6. Frequency table for primary schools with different student genders.**



**Figure 7. Distribution of primary schools with different student genders. (Upper: map with mark for all 3 categories; bottom: map with mark for only one-gender schools.)**

Almost all of the primary schools have both boys and girls (Co-Ed), while there are more Girl schools than Boy schools. The one-gender schools are located in only Hong Kong Island and Kowloon area.

**4.2.3 Religion**

```
Nonreligious                        238
PROTESTANTISM / CHRISTIANITY        183
CATHOLICISM                         107
BUDDHISM                             15
TAOISM                                8
CONFUCIANISM, BUDDHISM & TAOISM       4
ISLAM                                 2
OTHERS                                1
CONFUCIANISM                          1
Name: Religion, dtype: int64
```

**Figure 8. Frequency table for primary schools with different religions.**

Apart from the nonreligious schools (238 schools), the two most prominent groups of religious schools are Protestantism/Christianity (183 schools) and Catholicism (107) schools. There is not obvious pattern in its distribution on the map.



**Figure 9. Distribution of primary schools with different religions.**

**4.3 Evaluation of primary schools**

This score data for 100 primary schools is merged with the previous data for primary schools, with 94 overlapping items. Based on the data of these items, the

correlations between these factors were explored using "association" function from Dython package and a heatmap is generated (Figure 10).
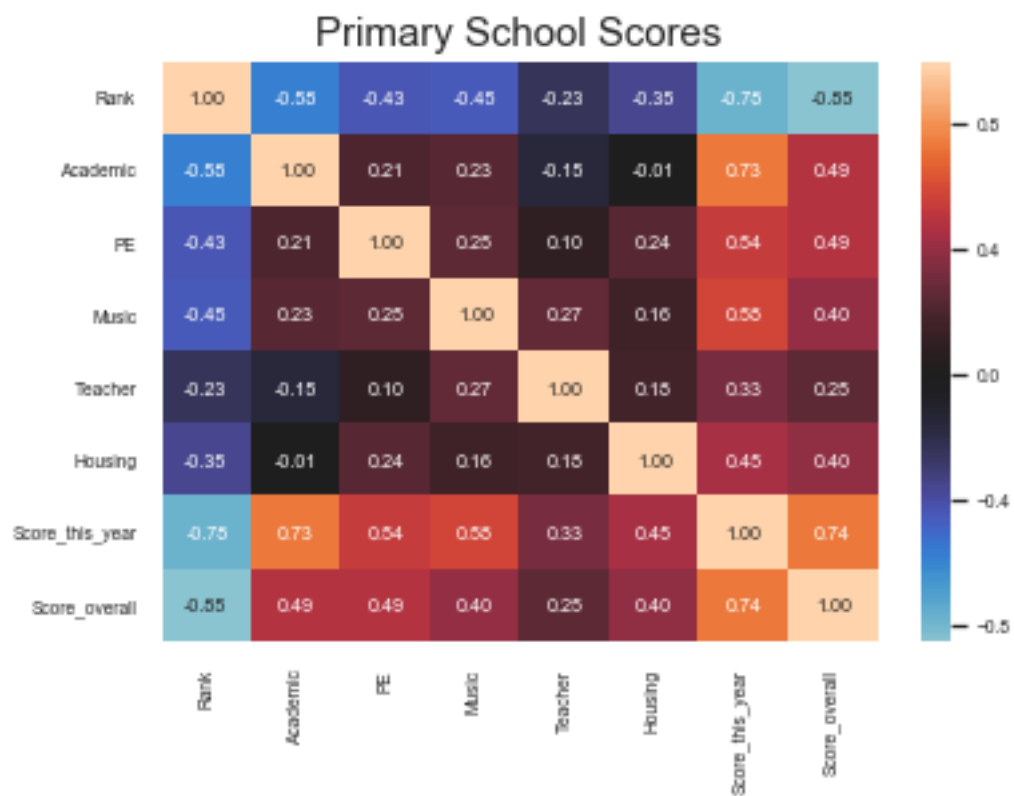


**Figure 10. Heatmap for correlations between different factors for school scores.**

As expected, ranking has a negative correlation with other factors, since the higher the score is, the smaller the rank number is. Also, it makes sense that the score this year has the highest correlation with all factors, since these are in principle derived from the other factors. The overall score also has a relatively high correlation with other factors, especially the score this year, indicating that the evaluation this year is largely consistent to the historical records. Surprisingly, if we leave aside the score this year and overall score, teacher has quite low correlation with all factors but music. On the other hand, housing is most highly correlated with PE. Especially, teacher and housing scores are even negatively correlated with academic. This probably indicates that music education is highly dependent on teachers, and PE is highly dependent on the place and environment, while academic is rather independent from these factors.

**4.4 Tuition of primary schools and correlation between different features of primary schools**

Tuition information were extracted for 54 primary schools, among which 53 have information for other features and only 15 also have information for scores. Based on the available information, the correlations between the previous analyzed features together with tuition and overall score were explored. The location represented by the district it belongs to is also included as a feature.
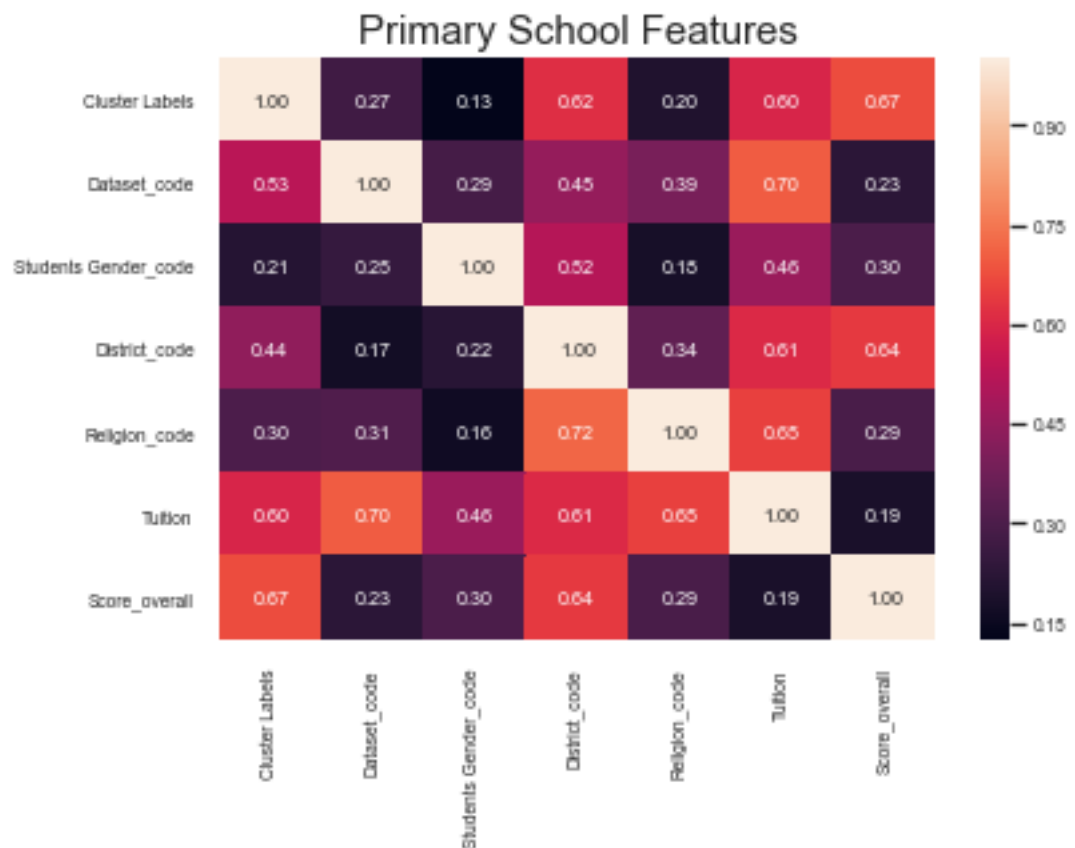


**Figure 11. Heatmap for correlations between different features for schools. For categorical-categorical correlations, row values indicate how well an index variable is predicted by column variables; column values indicate how well a column variable predicts row variables.**

Since the majority of the features is categorical, correlation ratio and Theil's U statistics were calculated, which are constantly positive. Among the features, tuition has a generally high correlation with other features except the overall score. This might be interpreted that high scored schools don't necessarily need higher tuition. However, since only 15 schools have information for both tuition and score, this value may not be very representative. Another interpretation is that, since these two factors are both continuous, the measurement of correlation is different from the other factors and therefore cannot be directly compared. Simply looking at the tuition and overall score, there is a weak positive correlation.
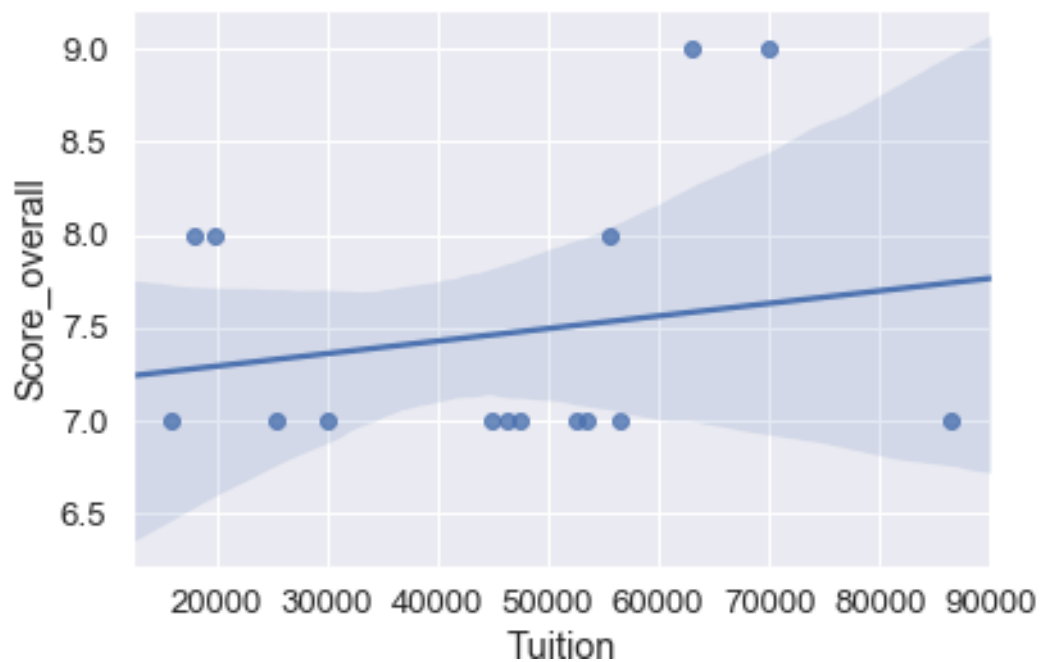
**Figure 12. Regression plot of tuition and overall score for 15 primary schools.**

Tuition has highest correlation with dataset, which makes sense since the dataset is defined by the funding type of the school, which should largely correlate with its tuition.
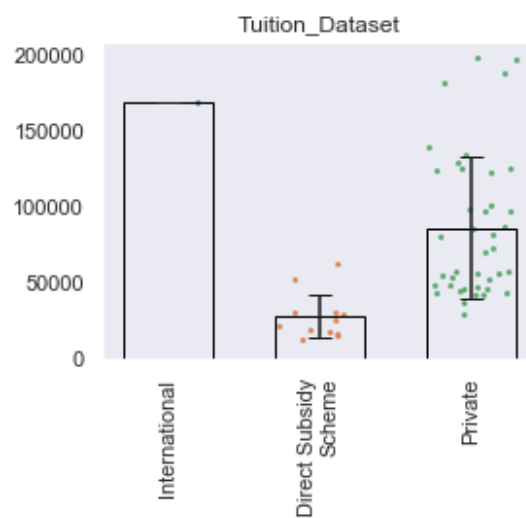


**Figure 13.   Tuitions for primary schools of different datasets.**

Only three types of schools are available for the tuition data and direct subsidy scheme schools have the lowest average tuition. Only one datapoint is present for international schools, making the average tuition the highest among the three. Private primary schools have a large range in tuition fee, which are generally higher than direct subsidy schools.
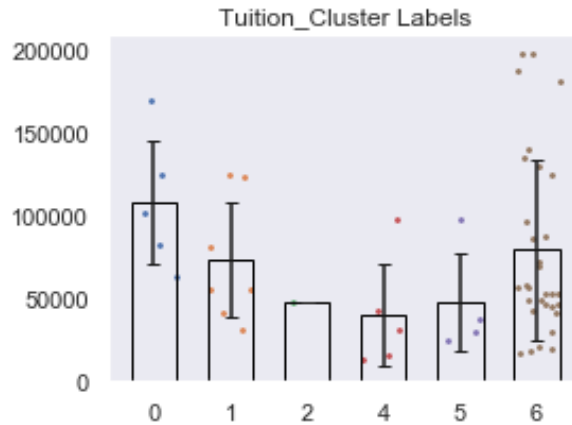
**Figure 14. Tuitions for primary schools of different cluster labels.**

Since cluster 6 contains majority of the schools, the tuition has a large variation in this cluster. Cluster 4 and 5 have relatively low tuition, which might be related to their locations away from downtown area.
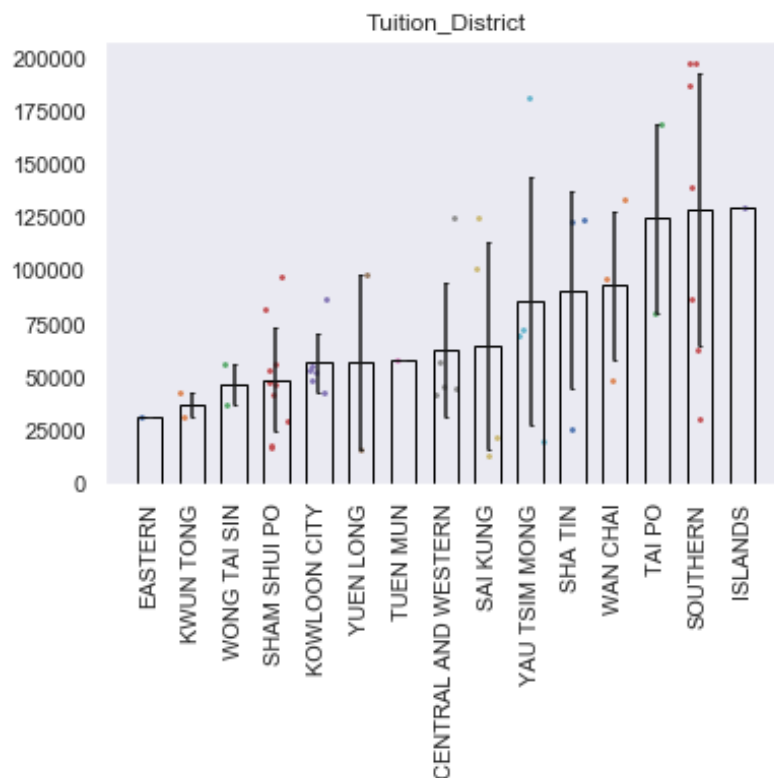


**Figure 15. Tuitions for primary schools of different districts.**

There was only one school in the Islands district, which represents the highest average tuition among districts. Eastern District also has one school only, which gives it the lowest average tuition. The small number of data points makes these two values less credible. Apart from them, Tai Po and Southern districts also

have high average tuition, while Tuen Muen, Wong Tai Sin, Sham Shui Po are on the lower end of the spectrum.
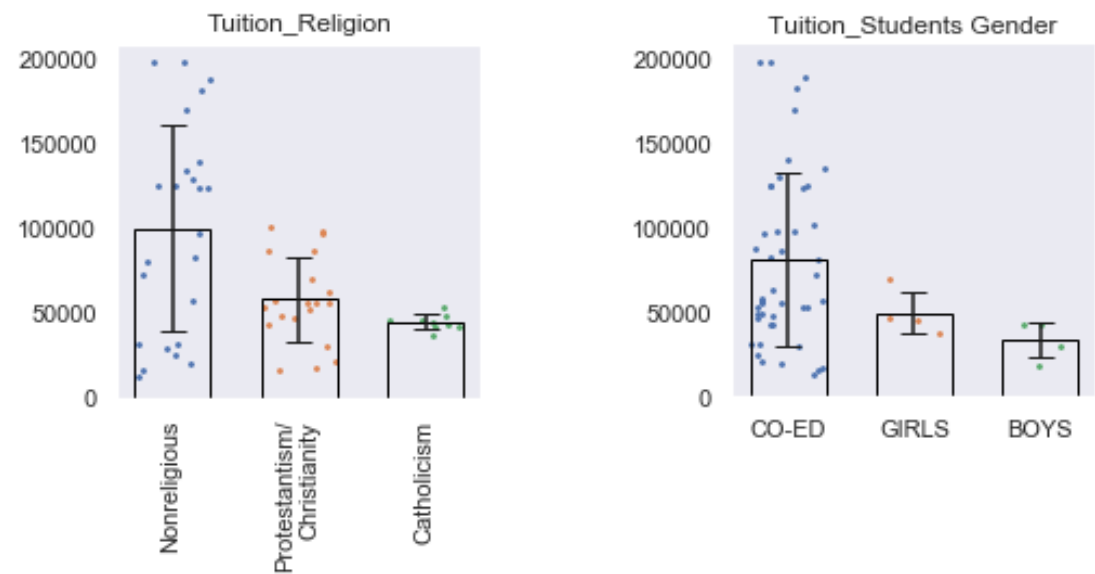


**Figure 16. Tuitions for primary schools of different religions and student genders.**

The comparison between schools of different religions and student genders were also explored. Generally, the most abundant group, which also represents the most unspecified group, has largest variation in tuition and higher average value as compared to the other more specified groups.

The heatmap also shows correlation between different features with the overall score, which is especially highly correlated with cluster label and district.
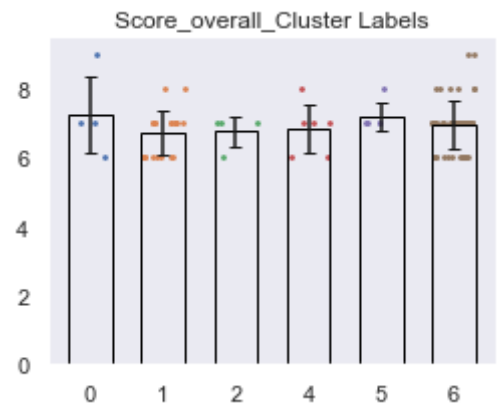


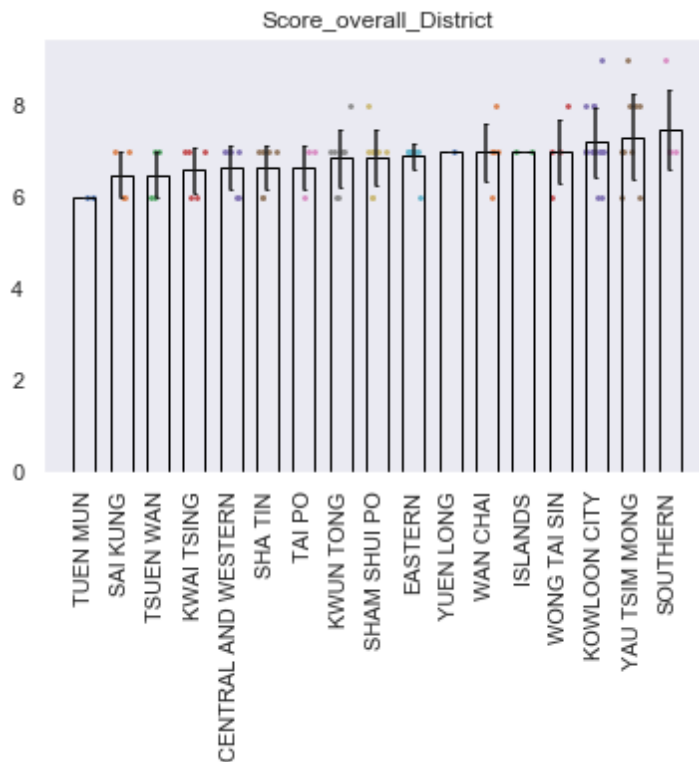**Figure 17. Overall scores for primary schools of different cluster labels.**

**Figure 18. Overall scores for primary schools of different districts.**

However, it could be seen that the average overall scores are not essentially different between the groups. However, a few schools with top scores are located in Kowloon City, Yau Tsim Mong and Southern districts, making them the three top districts in terms of school overall scores.
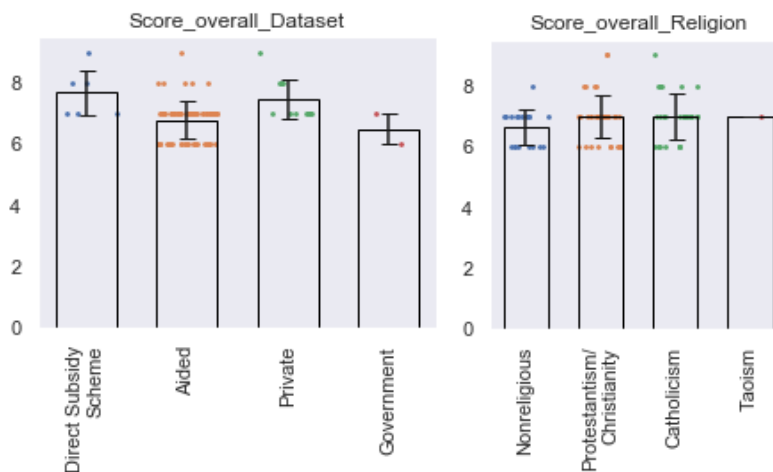


**Figure 19. Overall scores for primary schools of different datasets (left) and religions (right)**
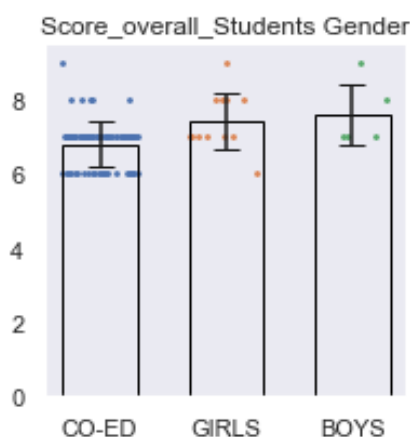
**Figure 20. Overall scores for primary schools of different student gender.**

The relationship between other factors with overall score were also explored. Although the difference is small, there is a trend for Direct subsidy scheme schools and private schools to have higher score as compared to aided and government schools; while single-gender schools tend to have higher score than Co-ed schools.

**4.5 Correlations between features of the 18 districts.**

Statistics about domestic income and rent was gathered for each district, including "Median Monthly Domestic Household Rent", "Median Monthly Domestic Household Income", "Median Rent to Income Ratio" and "Median Mortgage Payment and Loan Repayment to Income Ratio". The median household rend for different types of housing were also present in the dataset but were excluded due to the limitation of knowledge and difficulty in interpretation.

From the previous analysis, some information about primary schools in each district could be extracted, which could be correlated with other statistics of the district. The density of primary schools, average tuition and average overall score for each district, the relative fractions of schools from each cluster, and the relative fraction of schools from each dataset for each district were included in the analysis. Other features were excluded since they didn't show apparent difference between districts or are not as relevant. The correlations were analyzed using the same function and a heatmap was generated.
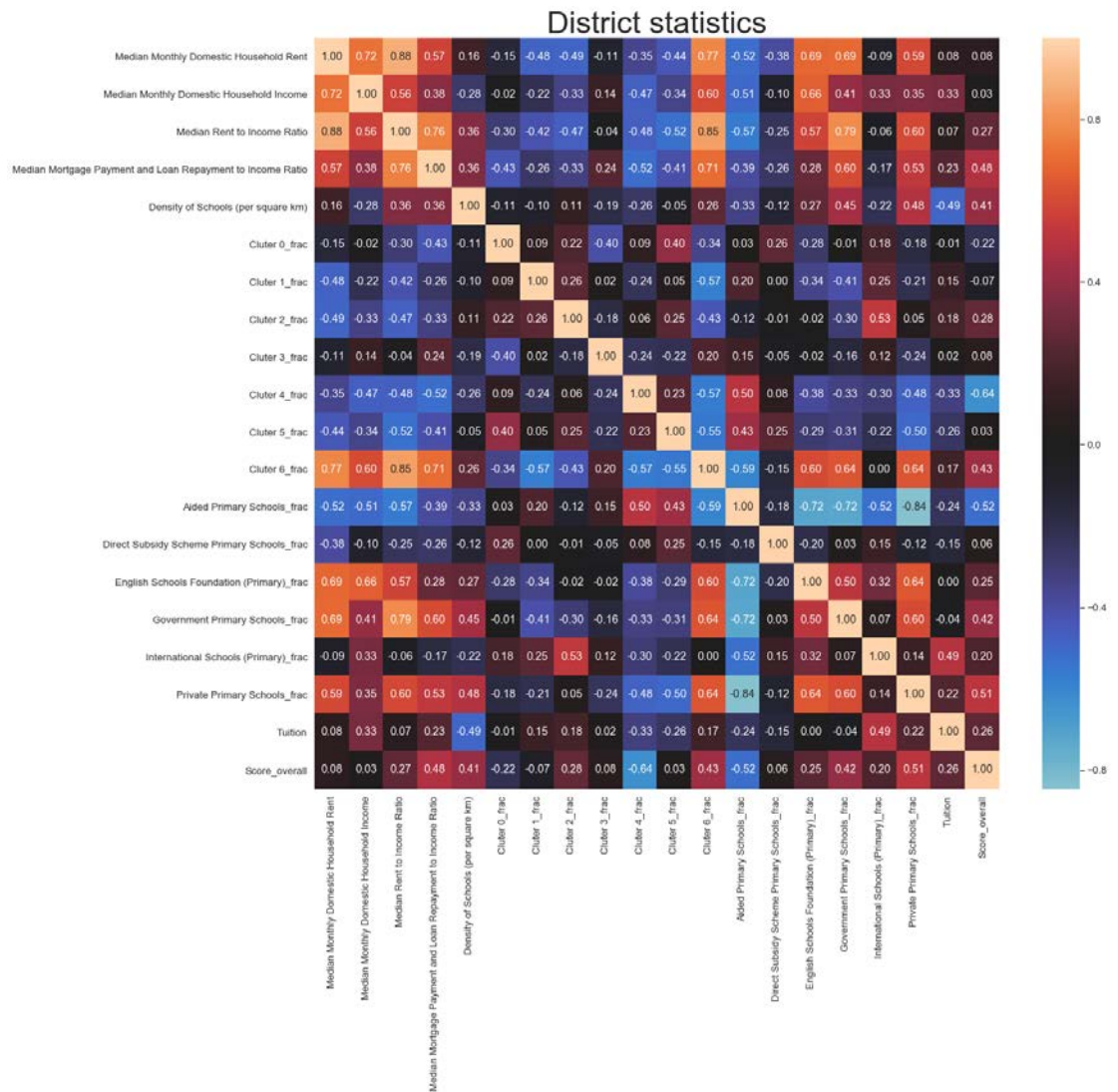
**Figure 21. Correlations between different features for each district.**

Since there are many variables and the relationships can be very complex, only a few observations would be discussed here and might be very incomplete.

One instant observation is that the four domestic income and rent statistics are positively correlated with each other. However, only the Median Monthly Domestic Household Income is negatively correlated with the density of schools.
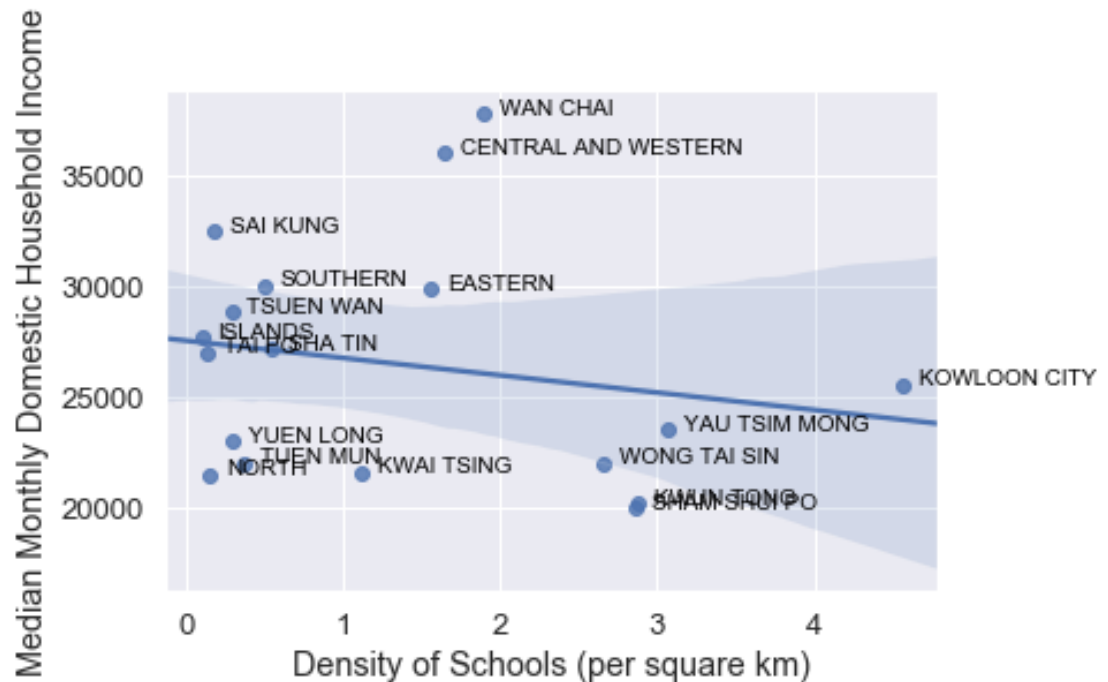
**Figure 22. Regression plot of Median Monthly Domestic Household Income and the density of schools in 18 districts.**

The regression plot shows that except from the two outliers (Wan Chai and Central and Western districts), the other districts generally align with the negative correlation between these two variables.

The second observation is that cluster 6 is the most distinct one as compared to other clusters, same as Aided Primary School in all dataset categories. This might be due to the large abundance of these two categories as compared to others, making them the dominant factor influencing the correlations of other categories.

To be specific, overall score is positively correlated with all dataset categories except Aided Primary Schools, which is consistent with previous analysis that Aided Primary Schools tend to have lower score. Average tuition has an apparent negative correlation with density of schools, which might be explained by the supply-demand relationship.

The pattern of correlations reveals a type of districts that is dominated by cluster 6 schools, which is correlated with lower proportion of Aided Primary Schools ($R^2$ = -0.59, Figure 23). This is consistent with the previous result that the dataset type can be largely predicted by the cluster label (Figure 11). These districts may also have higher density of schools, high "Median Monthly Domestic Household Rent", "Median Monthly Domestic Household Income", "Median Rent to Income Ratio" and "Median Mortgage Payment and Loan Repayment to Income Ratio",

given the positive correlations between these variables. Indeed, the regression plots between cluster 6 fraction and these variables (not shown) revealed Wan Chai, Central and Western, as well as Yau Tsim Mong as the three districts that are generally consistent with the previously described features, except that Yau Tsim Mong has lower income than many other districts.
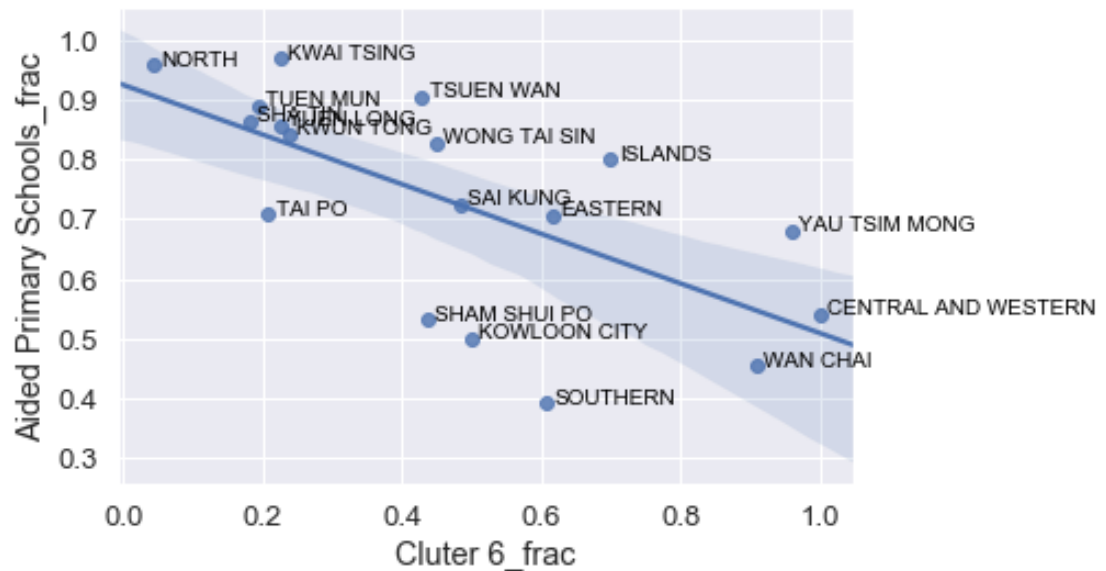


**Figure 23. Regression plot of fraction of cluster 6 schools and the fraction of aided primary schools.**

## 5. Discussion

There are interesting observations for the evaluation of schools, e.g., the score for teachers and housing is negatively correlated with the overall score, although the interpretation might be hard given very limited information. Also, higher score doesn't necessarily mean high tuition, although this correlation is only based on a small number of samples. Besides, direct subsidy scheme primary schools, private schools, and single-gender schools tend to have higher score than other types.

There are certainly many limitations for this project. One is that the features used for clustering is essentially not very differentiating and may not be directly related to other features (e.g. tuition and score, etc.). Also, it is possible to extract rating information for these venues for more interesting analysis that might be useful for business purposes, but was out of scope of my analysis here. Besides, the data for tuition and scores was from a single source that was not evaluated for the credibility due to the limit of resource.

## 6. Conclusion

This project attempts to have a broad exploration on the primary schools in Hong Kong in terms of their general features, geographical distribution, and nearby venues. The results revealed a largely similar pool of common venue categories near primary schools, which are mainly restaurants and food shops. However, a certain cluster (cluster 6) is predominant and seems to have a correlation with the most common type of school (aided primary school). The schools in this cluster is especially abundant in Wan Chai, Central and Western, and Yau Tsim Mong districts, which share many similarities in statistics like "Median Rent to Income Ratio". This indicates that the economic features of a district are correlated with the neighborhood and environment, including the localization of primary schools.

Overall, this project simply provides an overview about the general features and distributions of primary schools in Hong Kong. It may serve as a reference for parents to have an idea about the basic features of each district and the education resources. However, it is still far from extrapolating any social inferences or decisive conclusions from this simple analysis. A huge amount of additional information and research is surely required to reach more solid and comprehensive conclusions.

| Cluster Labels | 1st | 2nd | 3rd | 4th | 5th | 6th | 7th | 8th | 9th | 10th |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Chinese Restaurant | Fast Food Restaurant | Hong Kong Restaurant | Shopping Mall | Asian Restaurant | French Restaurant | Food Truck | Fried Chicken Joint | Food Court | Food & Drink Shop |
| 1 | Fast Food Restaurant | Chinese Restaurant | Shopping Mall | Hong Kong Restaurant | Park | Cha Chaan Teng | Bus Station | Asian Restaurant | Convenience Store | Food Court |
| 2 | Food Court | Food & Drink Shop | Park | Food Truck | Flower Shop | French Restaurant | Chinese Restaurant | Fried Chicken Joint | Fujian Restaurant | Garden |
| 3 | Shopping Mall | Light Rail Station | Fast Food Restaurant | Fried Chicken Joint | Food & Drink Shop | Flower Shop | Food Truck | Food Court | Cha Chaan Teng | Hong Kong Restaurant |
| 4 | Fast Food Restaurant | Shopping Mall | Hong Kong Restaurant | Cha Chaan Teng | Bus Station | Light Rail Station | Market | Chinese Restaurant | Park | Food Court |
| 5 | Coffee Shop | Chinese Restaurant | Café | Food & Drink Shop | Japanese Restaurant | Food Court | Fast Food Restaurant | Food Truck | Flower Shop | French Restaurant |
| 6 | Chinese Restaurant | Coffee Shop | Café | Fast Food Restaurant | Dessert Shop | Cantonese Restaurant | Noodle House | Shopping Mall | Japanese Restaurant | Cha Chaan Teng |