

徐海阳 PB20000326

Report for AD2022 Lab3

实验要求

实验要求

□ 实验要求

- 给定一个数据集和预测目标，需要分析数据、统计以及抽取特征
- 数据分析、统计，如：
 - 单个特征的分布
 - 统计缺失值
 - 特征间的相关性
 - 推测特征的含义
 - 异常样本
 - 数据抽样...
- 特征抽取，如：
 - 特征的变换，如：str 转 int, log
 - 尝试组合特征
 - 特征子集选择
 - ...

□ 数据集-PISA2015（筛减版）

- 本次试验针对PICA2015中的学生调查问卷数据集
- 助教已经做了筛选，现包含西语地区的32130个学生、429个特征

CNTSCHID	Region	STNAME	SUBNAME	SEX	AGE	Grade	Grade_F1	Grade_F2	STPROVINA	STFEDERNA	STFIDMUNA	STFIDMUNA	STFIDMUNA	STFIDMUNA
0	ESP00001	72403	ESP00017	72400001	2	2	1	1	1	1	1	1	1	1
1	ESP00001	72404	ESP00017	72400002	2	2	1	1	1	1	1	1	1	1
2	ESP00001	72405	ESP00017	72400003	2	2	1	1	1	1	1	1	1	1
3	ESP00001	72406	ESP00017	72400004	2	2	1	1	1	1	1	1	1	1
4	ESP00001	72407	ESP00017	72400005	2	2	1	1	1	1	1	1	1	1
...
32128	ESP00000	72416	ESP00014	72400001	2	2	1	1	1	1	1	1	1	1
32129	ESP00000	72416	ESP00014	72400002	2	2	1	1	1	1	1	1	1	1
32127	ESP00000	72416	ESP00014	72400003	2	2	1	1	1	1	1	1	1	1
32126	ESP00000	72416	ESP00014	72400004	2	2	1	1	1	1	1	1	1	1
32128	ESP00000	72416	ESP00014	72400005	2	2	1	1	1	1	1	1	1	1

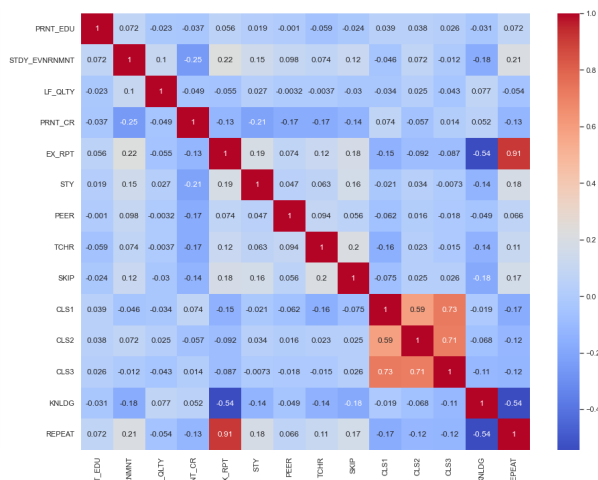
- CNTSCHID代表学校id，Region代表学生所在地区
- ST开头的是学生问卷回答特征
- 每个特征的具体含义可以参考codebook

□ Codebook是数据集每个特征的详细说明

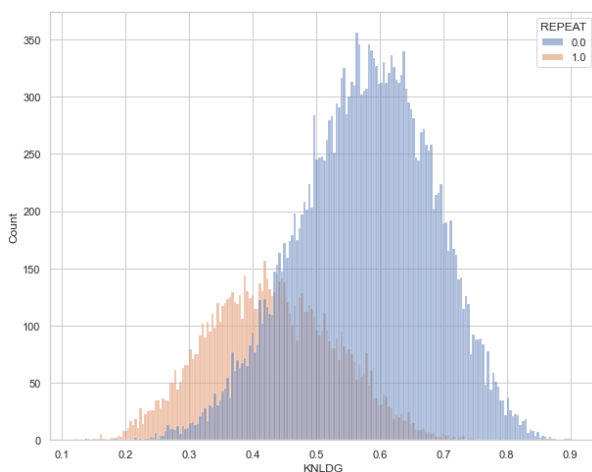
□ 预测任务

- 预测学生是否复读，即REPEAT列。

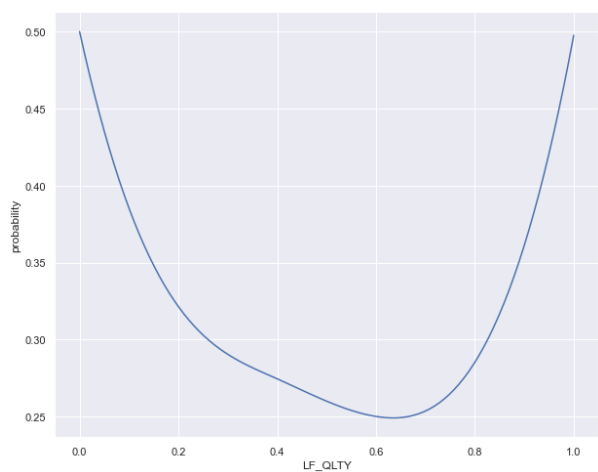
relationship between parental care and repetition



created new features heat map



relationship between knowledge richness and repetition



relationship between life quality and repetition

2 Approach

2.1 Preprocess

2.1.1 Data Cleaning

1. Data type uniform:

Turn str and object to int or float.

```
1 # temporarily turn nan to -1 for type conversion
2 df.replace(' ', -1, inplace=True)
3 # turn column[3] to str
4 df.iloc[:, 3] = df.iloc[:, 3].astype(str)
5 # turn the q part to `int`
6 df.iloc[:, 7:231] = df.iloc[:, 7:231].astype(int)
7 # turn the v part to `float`
8 df.iloc[:, 231:-1] = df.iloc[:, 231:-1].astype(float)
9 # actually null values are ' ' (space), need to change to NaN
10 df.replace(-1, np.nan, inplace=True)
```

2. Meaningless data discard:

Drop columns with the same data which means they are unrelated to the classification result.

```
1 # actually, we can see from the codebook that columns[5:19] are unrelated to the target variable
2 # so we can drop them
3 df.drop(columns=df.columns[5:19], inplace=True)
```

3. Normalization of abnormal data:

Pick out the data with vacant values, including the abnormal options in the questionnaire (9,99, 999, etc.)

And fill the data with the majority of that type of feature(mode number filling).

```
1 # process missing data which is xx9
2 # turn them into np.nan
3 max_list = []
4 for i in range(len(df.columns)):
5     ele = df.iloc[:, i].max()
6     max_list.append(ele)
7 for i in range(len(df.columns)):
8     if max_list[i] in [9,99,98,998,999,9.0,99.0,98.0,998.0,999.0,95.0]:
9         df.iloc[:, i].replace(max_list[i], np.nan, inplace=True)
10
11 # search all the column name with null values
12 null_columns = df.columns[df.isnull().any()]
13 # count the number of null values in each column, see those who have more than 10% null value
14 df[null_columns].isnull().sum()>3000
15 # actually those who have more than 10% null value are not meaningless, we can neglect them
16 # fill the null values with the mode of the column
17 mode_dict = df[null_columns].mode().iloc[0].astype(int).to_dict()
18 df.fillna(mode_dict, inplace=True)
```

4. Data standarlization:

MinMaxscaler to take all the data in [0,1].

```
1 # use MinMaxScaler to normalize the data
2 df_scaled_np = MinMaxScaler().fit_transform(df)
3 df_scaled = pd.DataFrame(df_scaled_np, columns=df.columns)
```

5. Auxiliary task:

Construct the relationship between paraphrases and numbers in the codebook and store them in the dictionary.

For subsequently easy view and operation.

```
1 # Below is to see the corresponding meaning of each columns and save to dicts,
2 # which will give great convenience when I want to refer the meaning of each feature.
3 cb = pd.read_excel('data/codebook.xlsx')
4 cb.NAME = cb.NAME.astype(str)
5 cb.drop(index = cb[cb.NAME=='nan'].index, inplace=True)
6 col = df.columns.to_numpy()
7 name = cb.NAME.to_numpy()
8 zai = []
9 for i in range(len(name)):
10     if name[i] in col:
11         zai.append(False)
12     else:
13         zai.append(True)
14 cb.drop(index = cb[zai].index, inplace=True)
15 # load cb's "NAME" and "VARLABEL" to a dict
16 cb_dict = cb.set_index('NAME').VARLABEL.to_dict()
17 # load cb_dict to "data/dictionary_cleaned.txt"
18 cb_dict_path = "data/dictionary_cleaned.txt"
```

```

19 | f = open(cb_dict_path, 'w')
20 | f.write(json.dumps(cb_dict, indent=0))

```

2.1.2 Feature Classification Manually

A number of representative feature classes were manually selected from over 400 features. The serial number represents the column number. Each feature class contains multiple columns of data.

parent education: 11~20 [PRNT_EDU](#)

study environment: 11 13~16 20 21 [STDY_EVNRMT](#)

life quality: 12 18 19 23-26 27-35 [LF_QLTY](#)

parent care: 36-39 [PRNT_CR](#)

former repeat: 48-50 [EX_RPT](#)

study attitude: 52-54 55 56 57-64 [STY](#)

peer relationship: 73-78 [PEER](#)

teacher relationship: 79-84 [TCHR](#)

skip lesson: 90-92 [SKIP](#)

class phenonmenon: 116-120, 121-129, 130-146 [CLS1, 2, 3](#)

knowledge: 304-403 [KNLDG](#)

2.1.3 Create Tool Function

```

1 | def make_df(df, index_list, add_repeat = False):
2 |     """
3 |     Select columns from `df` according to `index_list`, return a new df. If `add_repeat` == `True`, then the df will
4 |     automatically add the `REPEAT` column from the PISA data.
5 |     """
6 |     df_new = pd.DataFrame()
7 |     for i in range(len(index_list)):
8 |         index = index_list[i]
9 |         if (type(index) == int):
10 |             col = df.columns[index]
11 |         else:
12 |             col = index
13 |         df_new[col] = df[col]
14 |     if (add_repeat):
15 |         df_new['REPEAT'] = df['REPEAT']
16 |     return df_new
17 |
18 | def sum(df, name, mean = False):
19 |     """
20 |     Calculate the sum of all features for each sample in `df` simply by add one by one, return a new_df with only one column whose
21 |     name is `name`.
22 |     """
23 |     if (df.columns[-1] == 'REPEAT'):
24 |         df.drop(columns=['REPEAT'], inplace=True)
25 |     if (type(df) == pd.core.frame.DataFrame):
26 |         df = df.to_numpy()
27 |     long = len(df[0])
28 |     df = np.sum(df, axis=1)
29 |     if (mean):
30 |         df = df/long
31 |     df = pd.DataFrame(df, columns=[name])
32 |     return df
33 |
34 | def sum_stdiz(map, df, name):
35 |     """
36 |     Calculate the sum of all features for each sample in `df` according to the correlation of features revealed by the heatmap
37 |     `map`, return a new_df with only one column whose name is `name`.
38 |     """
39 |     if (type(df) == pd.core.frame.DataFrame):
40 |         df = df.to_numpy()
41 |     scaler = MinMaxScaler()
42 |     pn = map.iloc[-2]
43 |     sign = map.iloc[-1][0]
44 |     for i in range(len(pn)-1):
45 |         if (i==0):
46 |             if (sign>0):
47 |                 temp = df[:,i]
48 |             else:
49 |                 temp = -df[:,i]
50 |         else:
51 |             if ((pn[i]>0 and sign > 0) or (pn[i]<0 and sign < 0)):

```

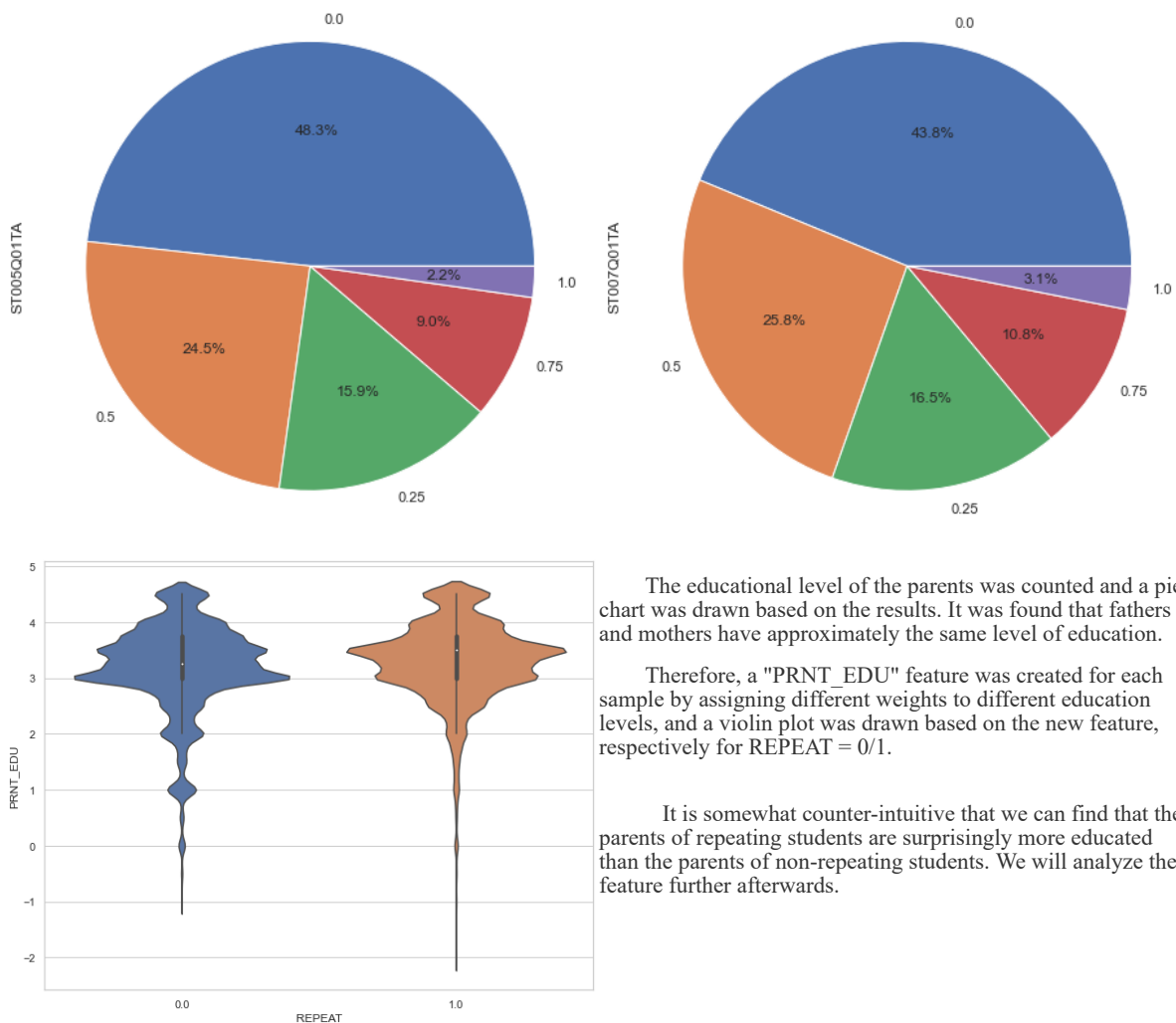
```

49         temp = temp + df[:,i]
50     else:
51         temp = temp - df[:,i]
52     df_new = scaler.fit_transform(temp.reshape(-1,1))
53     df_new_pd = pd.DataFrame(df_new).rename(columns={0: name})
54     return df_new_pd
55
56 def draw_corr_map(df,method='pearson'):
57     """
58     Calculate the covariance between all the characteristic variables in `df`, get the covariance matrix, and plot the heatmap
59     """
60     map = df.corr(method=method)
61     sns.heatmap(map, annot=True, cmap='coolwarm')
62     return map
63
64 # draw different plot
65 sns.violinplot
66 sns.histplot
67 sns.pieplot

```

2.2 Data Extraction and Analysis

2.2.1 PRNT_EDU(Parent Education)

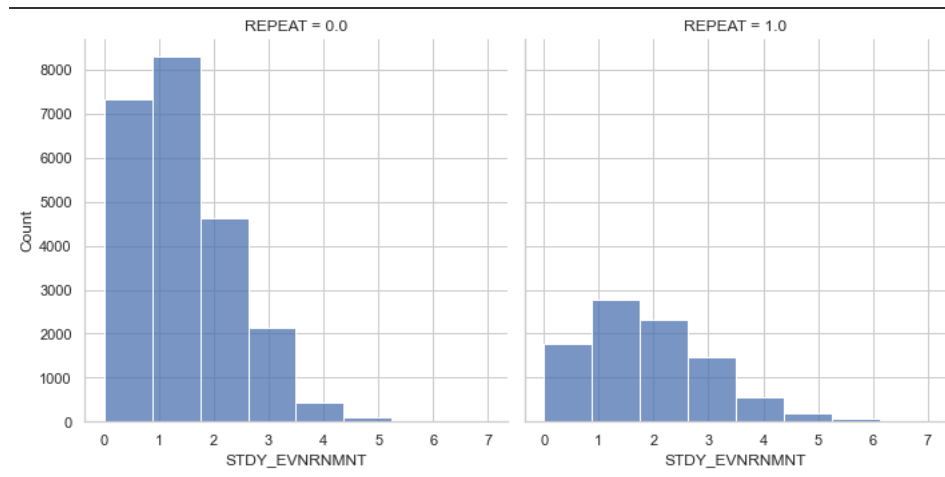


The educational level of the parents was counted and a pie chart was drawn based on the results. It was found that fathers and mothers have approximately the same level of education.

Therefore, a "PRNT_EDU" feature was created for each sample by assigning different weights to different education levels, and a violin plot was drawn based on the new feature, respectively for REPEAT = 0/1.

It is somewhat counter-intuitive that we can find that the parents of repeating students are surprisingly more educated than the parents of non-repeating students. We will analyze the feature further afterwards.

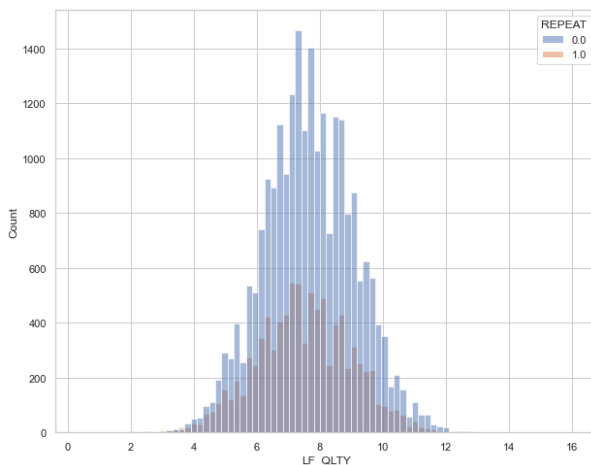
2.2.2 STDY_EVNRMNT(Study Environment)



Study environment is an indicator of the richness of electronic devices in the student learning environment. The "STDY_EVNRNMNT" characteristic is obtained by simply summing the contributions of various electronic devices. The frequency distribution of the feature is plotted separately for repeating or not repeating.

We can see that the richer the electronic devices, the higher the repetition rate of students. This is also in line with common sense, because with high electronic devices, students tend to be more likely to use them for things that are not related to their studies, such as playing games, chatting on social software, etc.

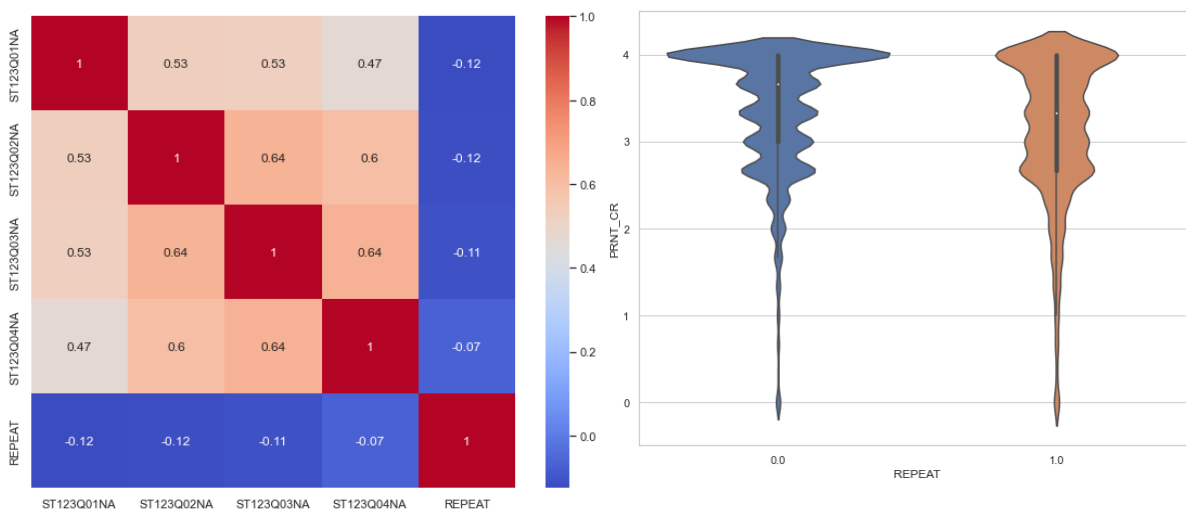
2.2.3 LF_QLTY(Life Quality)



"Life Quality" is relatively simple to obtain: the values of the variables related to household wealth and quality of life are simply summed. Histogram of the frequency distribution for repeating or not repeating.

We don't seem to see much difference in the distribution information from this graph, which we will discuss in detail later.

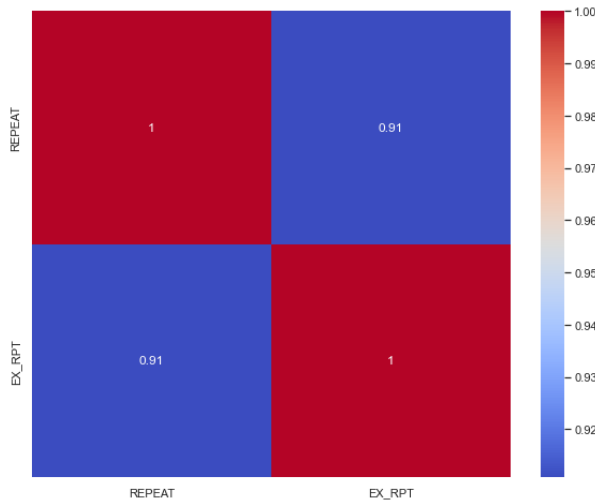
2.2.4 PRNT_CR(Parent Care)



A covariance matrix is calculated and a heat map is drawn based on the characteristics that show the degree of parental concern and encouragement for the student. Sum all the care up. Then violin plots were drawn based on whether to repeat and the parental care.

It can be found that these characteristics are highly correlated, which is in line with common sense, since parents who encourage students to get good grades tend to comfort them when they lose as well. Also, according to the violin chart we can see that students who do not repeat tend to have more parental care compared to those who do. And it's easy to understand too.

2.2.5 EX_RPT(Former REPEAT Record)

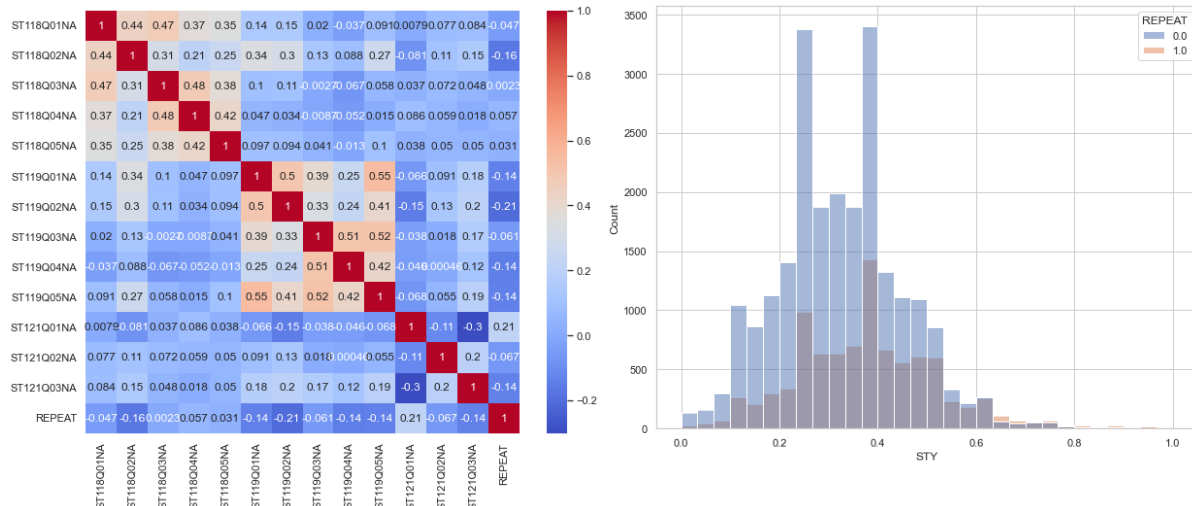


Past REPEAT records. When I was manually organizing the features, I realized that this was in fact a quickpath for predicting REPEAT. Let's see the data analysis.

Sum the three past repeat record together to form new feature "EX_RPT". Directly plot the relevance heat map based on the distribution. Good lord, the degree of correlation is almost linear! Perfect feature.

It's also understandable that a person often repeats because of various reasons, such as poor academic ability, too much negative emotion given by teachers and parents, and so on. Unless he makes changes, it is likely that he will continue to repeat.

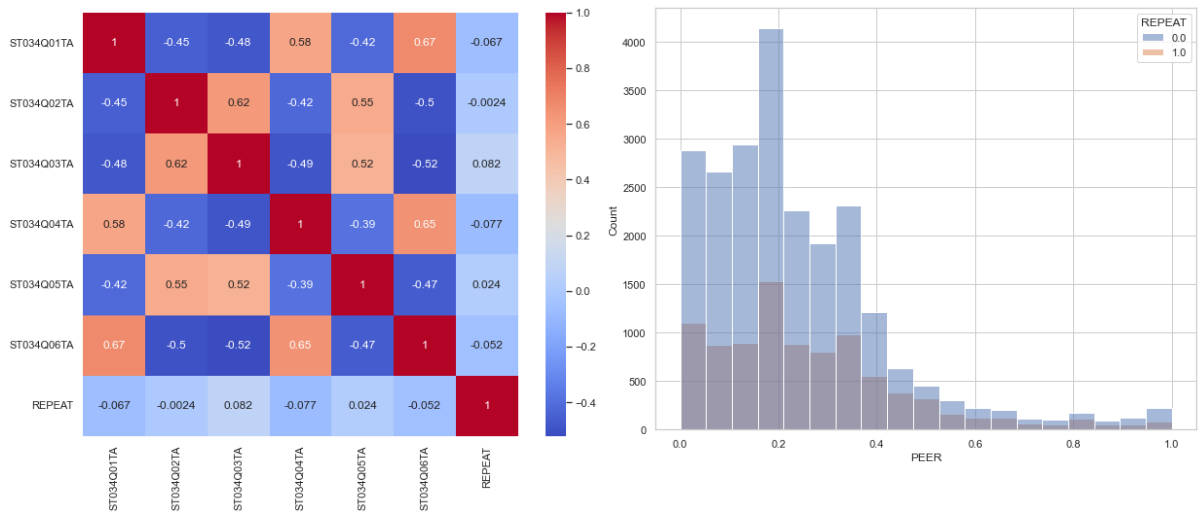
2.2.6 STY(STY Ambition)



Study ambition is certainly also an important factor influencing whether to repeat the study or not. List all the characteristics about study ambition and statistically plot the heat map. It is also weighted and normalized according to the map.

Then we draw the Histogram of the frequency distribution for repeating or not repeating. We can observe that students with lower ambition tend to have lower probability of repeating. I guess it is because of the high self-expectation, resulting in too much pressure, but not good learning.

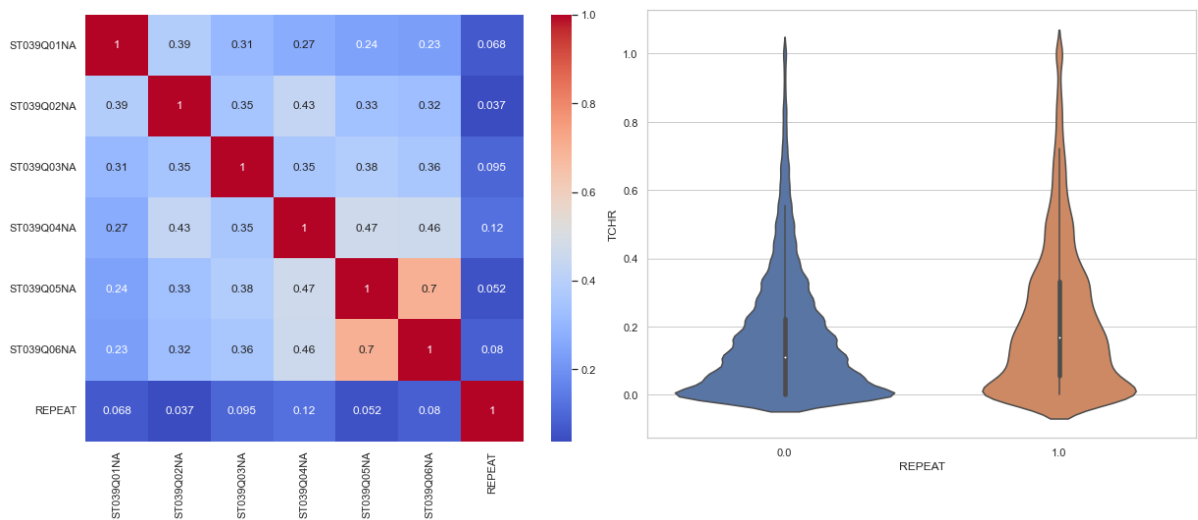
2.2.7 PEER(Peer Relationship)



Next is the PEER characteristic, which examines the relationship between partners. Heat maps are drawn based on statistical data and a strong correlation is found between the data. Therefore, the new feature "PEER" is obtained by fusing the features according to the map.

It can be found that compared to the normal repetition rate of 25%, the ratio of repetition to non-repetition among students with poor interpersonal relationships is about one to two, which is greater than 25%. This indicates that students who repeat tend to have poorer interpersonal relationships.

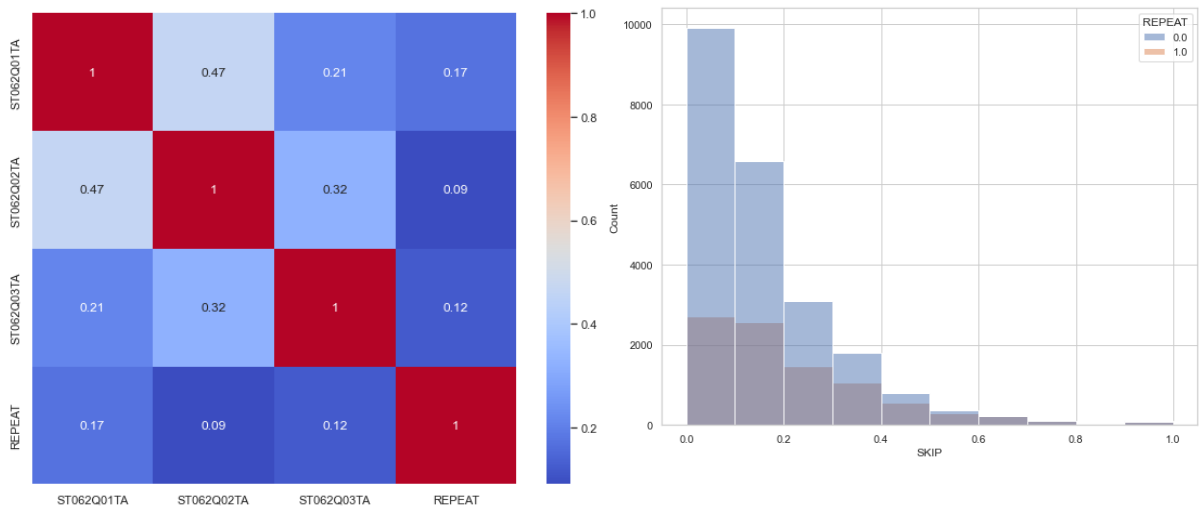
2.2.8 TCHR(Teacher Review)



TCHR characteristics are characteristics made up of your teacher's attitude and evaluation of you. The higher the value, the harsher and lower the teacher's evaluation of you. Draw a heat map and a violin map.

It can be found that students who repeat their studies tend to have a higher value of the TCHR characteristic, which indicates that the teacher's attitude and evaluation of them is worse and may even be offensive and insulting.

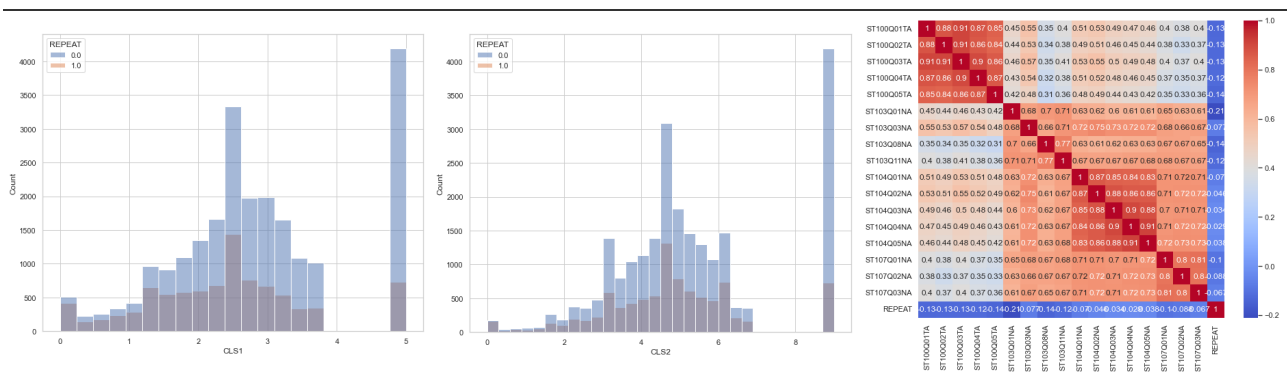
2.2.9 SKIP(Skip class)



SKIP characteristics mainly count the frequency and duration of students' skipping classes. The higher the value indicates the more serious students' skipping classes. Draw a heat map and a violin map.

We can find that students who repeat have a higher skipping rate, which also explains their worse academic performance in the opposite direction.

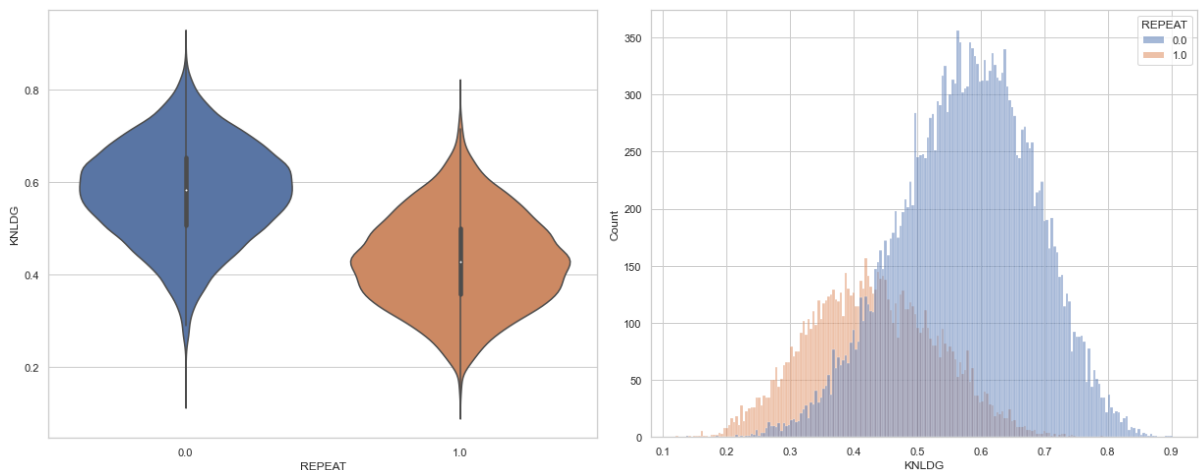
2.2.10 CLS(Class Teaching Phenomenon)



The CLS features counted a large number of data on classroom atmosphere and teachers' teaching ability, etc. We obtained CLS1, CLS2, and CLS3 features by summing up, which represent classroom teaching atmosphere, teacher's teaching ability, and course scheduling rationality, respectively. Frequency histograms of CLS1,2 and heat maps of CLS3 were drawn.

We can find that these manually selected features do possess an extremely strong correlation and can be divided into groups.

2.2.11 KNLDG(Knowledge)

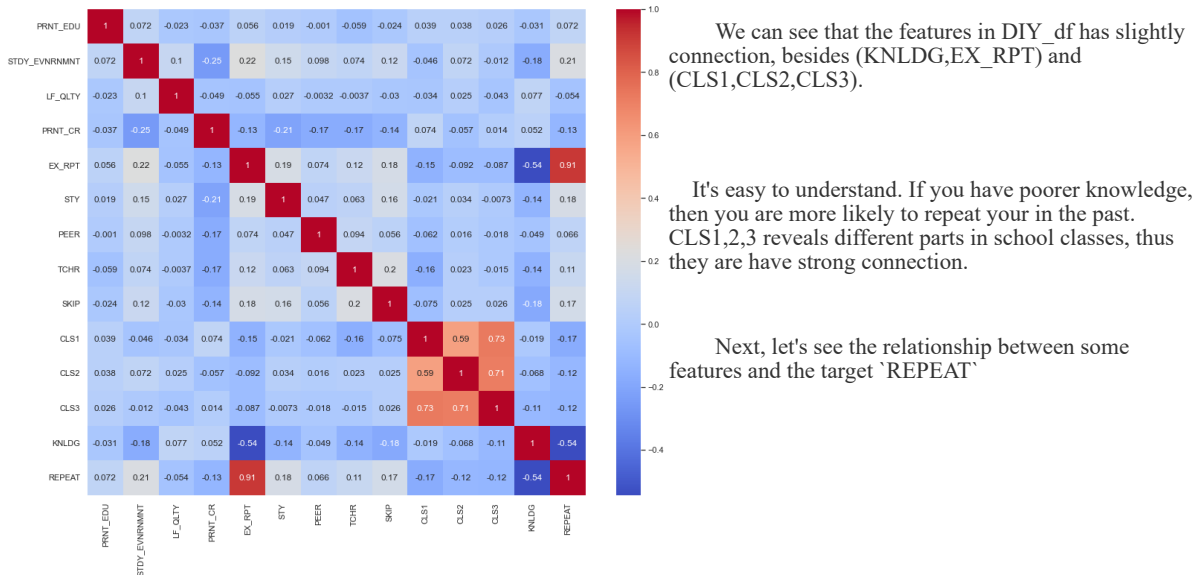


The Knowledge feature is also a quantity that I am sure must be highly correlated with REPEAT after constructing it.

Knowledge is a feature consisting of 100 academic test scores in a variety of scientific disciplines. Plotting frequency histograms and violin plots, we can see that students who do not repeat have a much higher level of knowledge compared to those who do. Perfect feature!

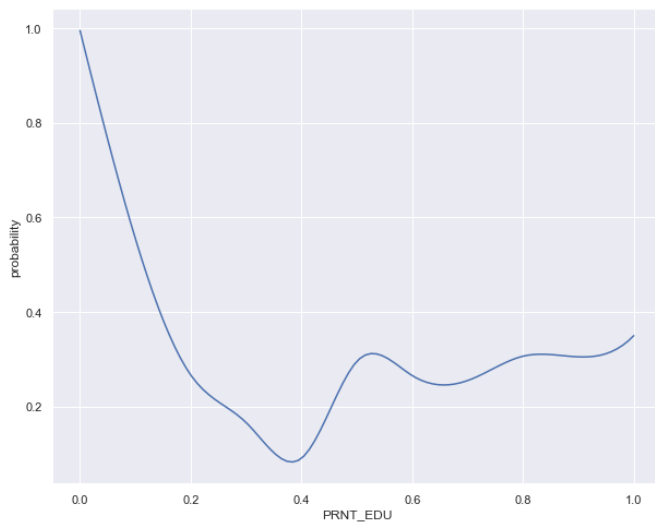
3 Further Analysis

Integrate all self-constructed features into a Dataframe. Count the covariance matrix and plot the heat map.



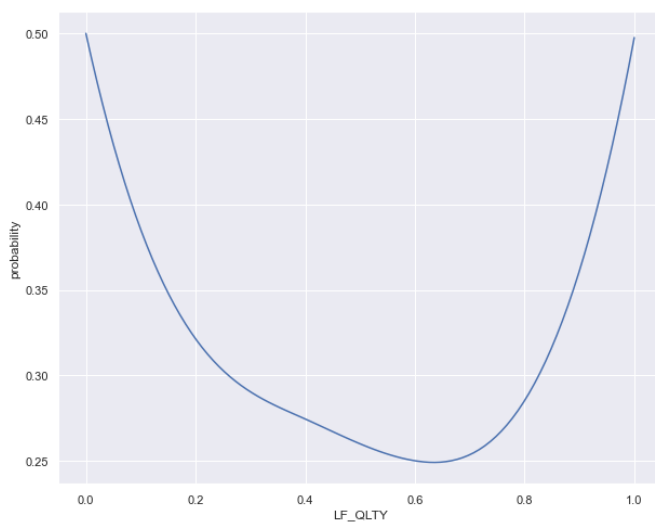
The data were discretized and the ratio between the number of newly increased repeaters in each small interval and the total number of repeaters and non-repeaters in that interval was counted, which is the repeater rate in that small interval, and from that the relationship curve between the repeater rate and the characteristic take values was drawn. Smoothing was performed with quadratic interpolation.

```
1 # discretize the data
2 def discrete(df, cond_col, step, tgt_col = 'REPEAT'):
3     # min, max = df[cond_col].describe()[3], df[cond_col].describe()[7]
4     sum = []
5     cur_sum = []
6     group = int(1/step)
7     for point in np.linspace(0, 1, group+1):
8         RPT, total = stat(df, cond_col, point, tgt_col, step)
9         cur_sum.append(RPT)
10        sum.append(total)
11    cur_sum = np.array(cur_sum)
12    sum = np.array(sum)
13    return cur_sum, sum
14
15 # count the total number in the group and the number of repeat in the group
16 def stat(df, cond_col, base, tgt_col, dur):
17     sum = 0.01
18     cur_sum = 0
19     index = df.index[(base <= df[cond_col]) & (df[cond_col] < base + dur)]
20     sub_df = pd.concat([df[cond_col][index], df[tgt_col][index]], axis=1)
21     for i in range(sub_df.shape[0]):
22         sum += 1
23         if sub_df[tgt_col].iloc[i] == 1:
24             cur_sum += 1
25     return cur_sum, sum
26
27 # draw the plot
28 def plt_discrete(df, cond_col, step, tgt_col = 'REPEAT'):
29     cur_sum, sum = discrete(df, cond_col, step, tgt_col)
30     x = np.linspace(0, 1, int(1/step)+1)
31     y = cur_sum/sum
32     itp = interpolate.CubicSpline(x, y) # smooth
33     xnew = np.linspace(0, 1, 100)
34     ynew = itp(xnew)
35     plt.plot(xnew, ynew)
36     plt.xlabel(cond_col)
37     plt.ylabel('probability')
38     plt.show()
```



First, the relationship between parents' education level and repetition rate is observed. It can be seen that children's repetition rate is lowest when parents' education level is around 0.4. Too little parental education leads to a significantly higher repetition rate than for those with relatively more educated parents, while parental education starts to affect the repetition rate less from about 0.5 and the curve starts to fluctuate.

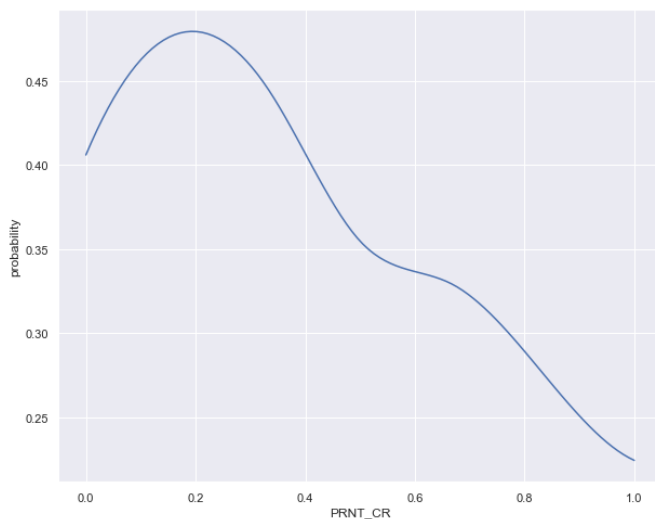
A wild guess is that parents with too little education will cause their children to grow up without much interest in learning or being able to supervise them. And after the parents' education level reaches a certain value, the level of attention to their children's learning will be about the same, and the impact will not be reflected in the repetition rate.



Next is the relationship between the quality of life level and the repetition rate. It can be significantly seen that the quality of life has the lowest repetition rate around 0.6. The repetition rate rises significantly at both ends of the spectrum.

It is well understood that if the quality of life is too low, students can hardly have the same study conditions to study; while if the quality of life is too high, students do not have the mind to study well, they don't have life pressure.

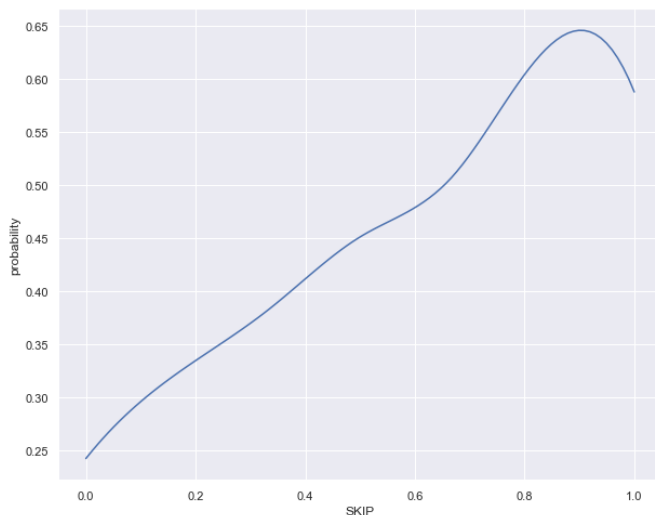
As a result, it is often the middle class students who will study hard, seeking to make the class leap and realize their own life values.



Next is the relationship between parental care and repetition rates. We can see that the repetition rate is highest when parental care is at 0.2. As parental care continues to increase, the repetition rate decreases significantly.

This is understandable. This is because the more encouragement and comfort from parents, the more the students' academic confidence grows and the relative stress is certainly not as great as that of students who are scolded by their parents every day.

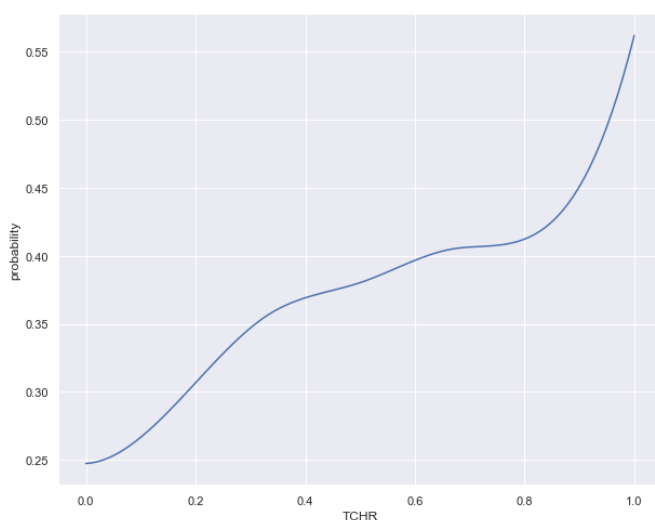
The reason why the repetition rate is slightly lower in places where parental care is close to 0 is, I think, probably because these students are self-motivated and indefatigable.



Next, we analyze the relationship between the skipping rate and the repeating rate. The repetition rate increases as the truancy rate increases, peaking at about 0.9.

This is easy to understand because the more you skip classes, the more likely it is that you are not learning enough. Therefore the repetition rate increases.

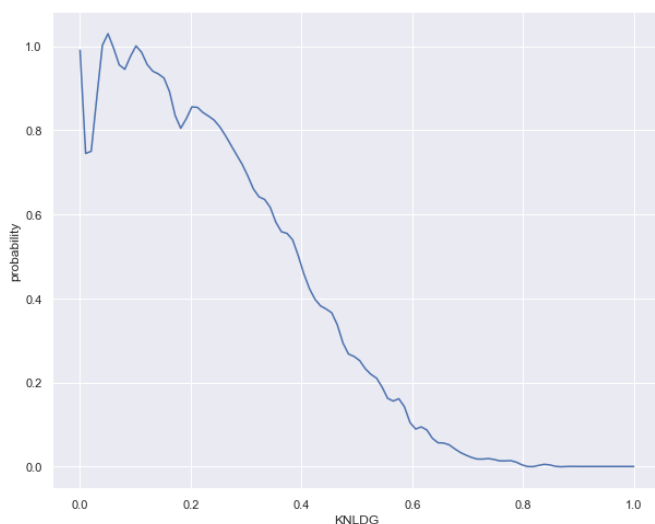
The reason why repetition rate decreases after 0.9, I think, is that some gifted students do not need the guidance of teachers, so the more they skip class, the more their academic performance goes up.



Next, we analyze the relationship between teachers' evaluations and students' repetition rates.

It can be seen that the harsher the teacher and the lower the evaluation of the student, the higher the repetition rate keeps rising.

This is an interactive process. Students often repeat because of poor academic performance, and teachers tend to look down on them or even insult them, which may further increase students' aversion to school and lead to a further decline in academic performance. A vicious circle.



Finally, there is the relationship between the student's knowledge base and the repetition rate. This is one of the clearest relationships.

The richer the knowledge base, the better the student's academic performance, and the less likely he or she is to repeat.

4 Supplementary

4.1 Environment and Packages

```
1 | import pandas as pd
2 | import matplotlib
3 | import matplotlib.pyplot as plt
4 | import seaborn as sns
5 | import numpy as np
6 | import json
7 | from sklearn.preprocessing import StandardScaler
8 | from sklearn.preprocessing import MinMaxScaler
9 | from scipy import interpolate
10 | %matplotlib inline
```

vscode

jupyter notebook

anaconda

4.2 Knowledge Covered

Data Preprocess :

1. Data Cleaning
2. Data Integration
3. Data Transformation
4. Data Statute

Feature Engineering :

1. Feature Extraction
2. Feature Design
3. Feature Construction

Covers Prof. QiLiu's AD2022 PPT slides Chapter2 : [Data Preprocess](#) and [Feature Engineer](#).

5 Conclusion

Lab3 costs me much more than 8h... About 3*8h... Sigh...

But I have learned a lot:

1. Feature Engineering
2. Convenience of Jupyter for showing and plotting
3. Seaborn Plotting of different kinds of excellent figures
4. CSS style design for prettier Markdown PDF