# China's K12 Online Education Data Management System

--based on Cassandra

# Part 1:

# Project Plan

## Topic:

China's K12 Online Education Data Management System

## Data Model:

Columnar

## Target Platform:

Cassandra (Datastax or Cosmos DB)

## Objective/Scope:

1.Study China's online education environment

2.Compare provided courses on different platforms(Such as  EDU,Zuoyebang,TAL,Yuanfudao etc)

3.Study the statistics of people using online education

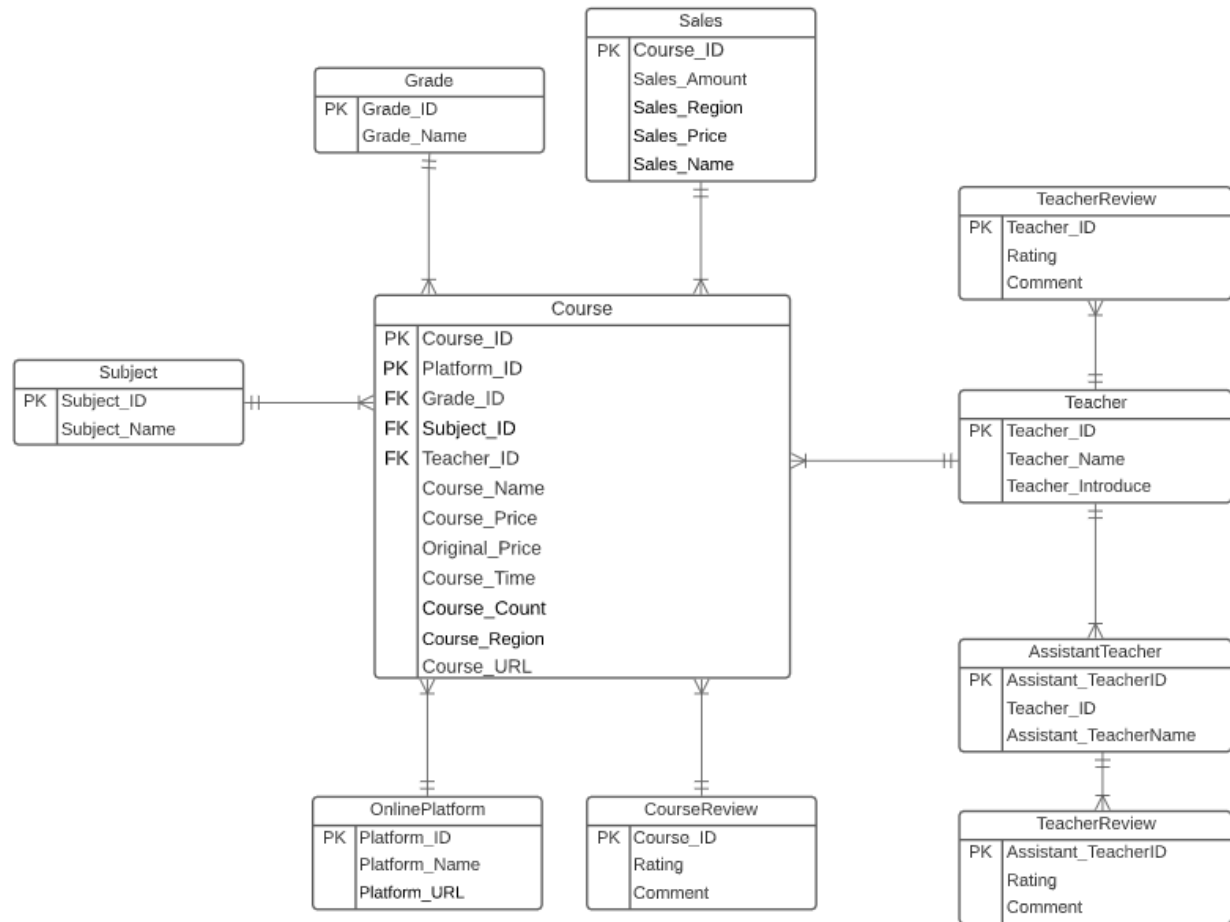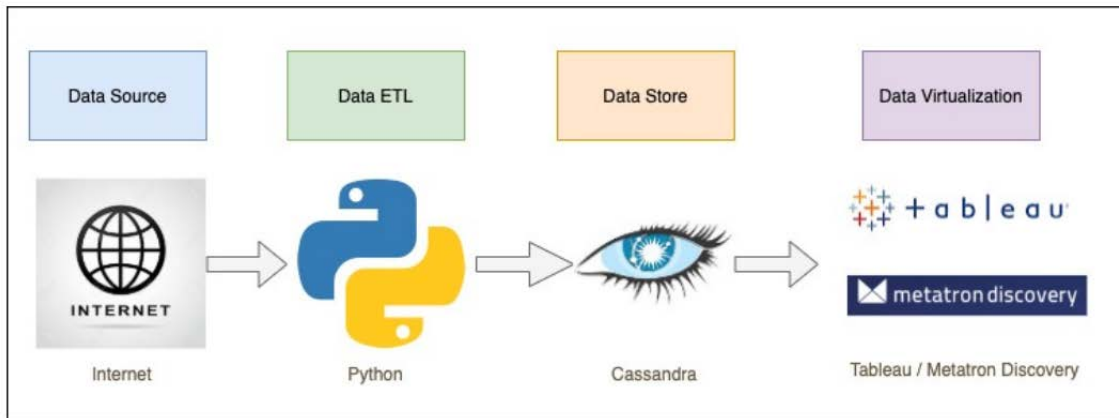4.Analyze growth rate in different platforms

## Visualizations Tool:

Tableau

# Part 2:

# Design

# ERD

# Architecture Diagram

# Part 3:

# Data Implementation

### 1) Data Source：

https://www.yuanfudao.com/



### 2) Web Scraping with Python:

```python
import requests
import json
import xlwt
import pandas as pd

listxueduan=['xiaoxue','chuzhong','gaozhong']
listgrade=['1','2','3','4','5','6','7','8','9','10','11','12']
list_c_1=['2','3','201']
list_c_2=['1','2','3','4','5','6','7','8','14']
list_c_3=['1','2','3','4','5','6','7','8','9']

lesson = []
teacher = []


def get_key_value(info_dict,k,j):
    try:
        info_dict[k] = j[k]
    except KeyError as e:
        for j_key in j:
            if isinstance(j[j_key], dict):
                get_key_value(info_dict,k,j[j_key])
list=[]
list11=[]
list12=[]
list13=[]


urllist=[]
```

```python
for i in listxueduan:
    if i=='xiaoxue':
        for j1 in range(1,7):
            #print(j1)
            g=j1
            for c1 in list_c_1:
                c=c1
                url = 'https://www.yuanfudao.com/tutor-student-
lesson/api/homepage?_productId=374&platform=www&version=5.11.0&UDID=4d27a58f757db800e
339317f1c245223&timestamp=1532658211041&startCursor=0&limit=18&grade='+str(g)+'&chann
elId='+str(c)+'&studyPhase='+str(i)+'&withNextGrade=false'
                urllist.append(url)
                list11.append(i)
                list12.append(j1)
                list13.append(c1)

    elif i=='chuzhong':
        for j2 in range(7,10):
            #print(j2)
            g=j2
            for c2 in list_c_2:
                c=c2
                url = 'https://www.yuanfudao.com/tutor-student-
lesson/api/homepage?_productId=374&platform=www&version=5.11.0&UDID=4d27a58f757db800e
339317f1c245223&timestamp=1532658211041&startCursor=0&limit=18&grade='+str(g)+'&chann
elId='+str(c)+'&studyPhase='+str(i)+'&withNextGrade=false'
                urllist.append(url)
                list11.append(i)
                list12.append(j2)
                list13.append(c2)
    else :
        for j3 in range(10,13):
            #print(j3)
            g=j3
            for c3 in list_c_3:
                c=c3
                url = 'https://www.yuanfudao.com/tutor-student-
lesson/api/homepage?_productId=374&platform=www&version=5.11.0&UDID=4d27a58f757db800e
339317f1c245223&timestamp=1532658211041&startCursor=0&limit=18&grade='+str(g)+'&chann
elId='+str(c)+'&studyPhase='+str(i)+'&withNextGrade=false'
                urllist.append(url)
                list11.append(i)
                list12.append(j2)
                list13.append(c3)
study_dict={
        'studyPhase':list11,
        'grade':list12,
        'channelId':list13
        }
def main():
    for i in range(0,72):
        url = urllist[i]
        s1=study_dict['studyPhase'][i]
        s2=study_dict['grade'][i]
        s3=study_dict['channelId'][i]
        r = requests.get(url)
```

```
        json_r= r.json()['list']
        for j in json_r:
            info_dict={
                    'studyPhase':s1,
                    'grade':s2,
                    'channelId':s3,
                    'id': None,
                    'minPrice': None,
                    'maxPrice': None,
                    'name': None,
                    'soldCount': None,
                    'price':None,
                    'subName': None,
                    'teachers':None
                    }
            for k in info_dict:
                get_key_value(info_dict,k,j)

            print(info_dict)
            lesson.append(info_dict)

            for t in j['teachers']:
                teacher_dict={
                        'lessonid':info_dict['id'],
                        'id':t['id'],
                        'nickname':t['nickname'],
                        'avatar':t['avatar']
                            }
                teacher.append(teacher_dict)


        work=xlwt.Workbook()
        sheet1=work.add_sheet('sheet1',cell_overwrite_ok=True)

head=['studyPhase','grade','channelId','id','minPrice','maxPrice','name','soldCount',
'price','subName']
        y=0
        for item in head:
            sheet1.write(0,y,item)
            y+=1
        x=1

        for item in lesson:
            if isinstance(item,dict):
                for head_item in head:
                    if head_item in item.keys():
                        y=head.index(head_item)
                        sheet1.write(x,y,item[head_item])
            x+=1
        work.save('yuanfudao new2.xls')

        work2=xlwt.Workbook()
        sheet1=work2.add_sheet('sheet1',cell_overwrite_ok=True)
        head2=['lessonid','id','nickname','avatar']
```

```
        y=0
        for item in head2:
            sheet1.write(0,y,item)
            y+=1
        x=1

        for item in teacher:
            if isinstance(item,dict):
                for head_item in head2:
                    if head_item in item.keys():
                        y=head2.index(head_item)
                        sheet1.write(x,y,item[head_item])
            x+=1
        work2.save('yuanfudao teacher.xls')


if __name__ == '__main__':
    main()
```

### 3) Rawdata

Data from (2020/01/01-2020/12/31,every month)



### 4) After data cleaning
a) course info table

| grade_id | course_id | course_name | start_time | end_time | lesson_num | run_id | teacher_id | subject_id | platform_id |
|---|---|---|---|---|---|---|---|---|---|
| 7 | 3511191 | 初一语文题型技巧特训班 | 2020-03-07 | 2020-03-10 | 7 | 2020-02-15 | 46 | 8 | 1 |
| 7 | 3163273 | 初一语文春季系统班（周五18:00) | 2020-02-21 | 2020-06-19 | 36 | 2020-02-15 | 46 | 8 | 1 |
| 7 | 3163275 | 初一语文春季系统班（周六09:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 46 | 8 | 1 |
| 7 | 3163277 | 初一语文春季系统班（周六14:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 49 | 8 | 1 |
| 7 | 3163279 | 初一语文春季系统班（周六14:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 46 | 8 | 1 |
| 7 | 3163281 | 初一语文春季系统班（周六18:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 49 | 8 | 1 |
| 7 | 3163283 | 初一语文春季系统班（周六18:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 52 | 8 | 1 |
| 7 | 3163285 | 初一语文春季系统班（周日09:00) | 2020-02-23 | 2020-06-21 | 36 | 2020-02-15 | 53 | 8 | 1 |
| 7 | 3163287 | 初一语文春季系统班（周日09:00) | 2020-02-23 | 2020-06-21 | 36 | 2020-02-15 | 52 | 8 | 1 |
| 7 | 3163289 | 初一语文春季系统班（周日18:00) | 2020-02-23 | 2020-06-21 | 36 | 2020-02-15 | 53 | 8 | 1 |
| 7 | 3511447 | 初一数学题型技巧特训班 | 2020-03-07 | 2020-03-10 | 7 | 2020-02-15 | 56 | 9 | 1 |
| 7 | 3168713 | 【目标A++班】初一数学春季系统班（周日14:00) | 2020-02-23 | 2020-06-21 | 36 | 2020-02-15 | 57 | 9 | 1 |
| 7 | 3168709 | 【目标A++班】初一数学春季系统班（周六18:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 57 | 9 | 1 |
| 7 | 3168543 | 【目标A+班】初一数学春季系统班（人教版·周五18:00) | 2020-02-21 | 2020-06-19 | 36 | 2020-02-15 | 59 | 9 | 1 |
| 7 | 3168545 | 【目标A+班】初一数学春季系统班（人教版·周五18:00) | 2020-02-21 | 2020-06-19 | 36 | 2020-02-15 | 56 | 9 | 1 |
| 7 | 3168547 | 【目标A+班】初一数学春季系统班（人教版·周五18:00) | 2020-02-21 | 2020-06-19 | 36 | 2020-02-15 | 61 | 9 | 1 |
| 7 | 3168549 | 【目标A+班】初一数学春季系统班（人教版·周五18:00) | 2020-02-21 | 2020-06-19 | 36 | 2020-02-15 | 62 | 9 | 1 |
| 7 | 3168551 | 【目标A+班】初一数学春季系统班（人教版·周六09:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 61 | 9 | 1 |
| 7 | 3168553 | 【目标A+班】初一数学春季系统班（人教版·周六09:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 64 | 9 | 1 |
| 7 | 3168555 | 【目标A+班】初一数学春季系统班（人教版·周六09:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 65 | 9 | 1 |
| 7 | 3168557 | 【目标A+班】初一数学春季系统班（人教版·周六14:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 66 | 9 | 1 |
| 7 | 3168559 | 【目标A+班】初一数学春季系统班（人教版·周六14:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 56 | 9 | 1 |
| 7 | 3168561 | 【目标A+班】初一数学春季系统班（人教版·周六14:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 68 | 9 | 1 |
| 7 | 3168563 | 【目标A+班】初一数学春季系统班（人教版·周六14:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 69 | 9 | 1 |
| 7 | 3168565 | 【目标A+班】初一数学春季系统班（人教版·周六18:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 61 | 9 | 1 |
| 7 | 3168567 | 【目标A+班】初一数学春季系统班（人教版·周六18:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 64 | 9 | 1 |
| 7 | 3168569 | 【目标A+班】初一数学春季系统班（人教版·周六18:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 72 | 9 | 1 |
| 7 | 3168571 | 【目标A+班】初一数学春季系统班（人教版·周六18:00) | 2020-02-22 | 2020-06-20 | 36 | 2020-02-15 | 58 | 9 | 1 |

b) sale info table:

| | course_id | price | original_price | quantity | sold_count | pre_sale_time | start_sale_time | end_sale_time | left_count | sold_out |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | | | | | | | | | |
| 2 | 3511191 | 9 | 299 | 80 | 7 | | | 2020-03-01 | 73 | 0 |
| 3 | 3163273 | 1200 | 0 | 30 | 26 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 4 | 0 |
| 4 | 3163275 | 1200 | 0 | 30 | 21 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 9 | 0 |
| 5 | 3163277 | 1200 | 0 | 30 | 26 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 4 | 0 |
| 6 | 3163279 | 1200 | 0 | 30 | 13 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 17 | 0 |
| 7 | 3163281 | 1200 | 0 | 30 | 22 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 8 | 0 |
| 8 | 3163283 | 1200 | 0 | 30 | 10 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 20 | 0 |
| 9 | 3163285 | 1200 | 0 | 30 | 22 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 8 | 0 |
| 10 | 3163287 | 1200 | 0 | 30 | 29 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 1 | 0 |
| 11 | 3163289 | 1200 | 0 | 30 | 8 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 22 | 0 |
| 12 | 3511447 | 9 | 299 | 80 | 12 | | | 2020-03-01 | 68 | 0 |
| 13 | 3168713 | 1200 | 0 | 30 | 24 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 6 | 0 |
| 14 | 3168709 | 1200 | 0 | 30 | 23 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 0 | 1 |
| 15 | 3168543 | 1200 | 0 | 30 | 16 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 14 | 0 |
| 16 | 3168545 | 1200 | 0 | 30 | 20 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 10 | 0 |
| 17 | 3168547 | 1200 | 0 | 30 | 21 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 9 | 0 |
| 18 | 3168549 | 1200 | 0 | 30 | 16 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 14 | 0 |
| 19 | 3168551 | 1200 | 0 | 30 | 16 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 14 | 0 |
| 20 | 3168553 | 1200 | 0 | 30 | 24 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 6 | 0 |
| 21 | 3168555 | 1200 | 0 | 30 | 23 | 2019-10-26 | 2019-10-26 | 2020-04-07 | 7 | 0 |

c) teacher table

| teacher_id | teacher_name |
|---|---|
| 0 | 刘薇 |
| 1 | 程磊 |
| 2 | 严攀 |
| 4 | 斯琴 |
| 5 | 满晓桐 |
| 6 | 苗锟 |
| 7 | 周梦麟 |
| 8 | 梁冰 |
| 9 | 平赫 |
| 11 | 铁健栩 |
| 12 | 李林 |
| 15 | 刘和妍 |

d) subject table

| subject_id | subject_name | subject_name_eng |
|---|---|---|
| 9 | 数学 | math |
| 1 | 政治 | politics |
| 2 | 地理 | geography |
| 3 | 历史 | history |
| 4 | 生物 | biology |
| 5 | 化学 | chemistry |
| 6 | 物理 | physics |
| 7 | 英语 | enghlish |
| 8 | 语文 | literature |

e) grade table

| grade_id | grade_name |
|---|---|
| 1 | 一年级 |
| 2 | 二年级 |
| 3 | 三年级 |
| 4 | 四年级 |
| 5 | 五年级 |
| 6 | 六年级 |
| 7 | 七年级 |
| 8 | 八年级 |
| 9 | 九年级 |
| 10 | 高一 |
| 11 | 高二 |
| 12 | 高三 |

## 5) Create table and Insert data

1. Connect to the Cassandra on datastax by using cassandra-driver



2.Create table



3.Insert data



5) Check data in datastax

**CQL console:**

Platform table:

```
token@cqlsh:info7275> select * from platform;

 platform_id | platform_name | platform_url
-------------+---------------+---------------------------
           3 |        youdao |       https://ke.youdao.com/
           2 |    genshuixue | https://www.genshuixue.com/
           1 |     yuanfudao |  https://www.yuanfudao.com/

(3 rows)
```

Subject table:

```
token@cqlsh:info7275> select * from subject;

 subject_id | subject_name | subject_name_eng
------------+--------------+-------------------
          6 |         物理 |          physics
          7 |         英语 |         enghlish
          9 |         数学 |             math
          4 |         生物 |          biology
          3 |         历史 |          history
          5 |         化学 |        chemistry
          8 |         语文 |       literature
          2 |         地理 |        geography
          1 |         政治 |         politics
```

Grade table:

```
token@cqlsh:info7275> select * from grade;

 grade_id | grade_name
----------+------------
        6 |      六年级
        7 |      七年级
        9 |      九年级
       10 |        高一
        4 |      四年级
        3 |      三年级
        5 |      五年级
        8 |      八年级
        2 |      二年级
       12 |        高三
       11 |        高二
        1 |      一年级
```

Teacher table:



Course table:



Sales table：

# Part 4:

# Data Visualization

# China's K12 Online Education Data Management System

## Course Count by Grade



| | Kindergart.. | Grade 1 | Grade 2 | Grade 3 | Grade 4 | Grade 5 | Grade 6 | Grade 7 | Grade 8 | Grade 9 | Grade 10 | Grade 11 | Grade 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 335 | 508 | 536 | 562 | 525 | 474 | 477 | 514 | 485 | 470 | 572 | 134 | 195 |

Line chart (Course Count): Feb 2020: 525, Mar 2020: 78, Apr 2020: 49, May 2020: 74, Jun 2020: 38, Jul 2020: 317, Aug 2020: 754, Sep 2020: 788, Dec 2020: 2, Jan 2021: 168, Feb 2021: 351, Mar 2021: 360

## Course Count Percent By Subject

Pie chart: 1.0%, 4.2%, 16.8%, 1.2%, 19.0%, 48.3%

## Course Count By Subject

| Subject | Count |
|---|---|
| Politics | 45 |
| Geography | 54 |
| History | 55 |
| Biology | 77 |
| Chemistry | 197 |
| Physics | 322 |

# China's K12 Online Education Data Management System

## Course Average Price by Grade



| Avg. Price (RMB) | | |
|---|---|---|
| Kindergarten | 1,005 | |
| Grade 1 | 975 | |
| Grade 2 | 919 | |
| Grade 3 | 937 | |
| Grade 4 | 917 | |
| Grade 6 | 932 | |
| Grade 7 | 1,012 | |
| Grade 8 | 1,014 | |
| Grade 9 | 1,079 | |
| Grade 10 | 1,063 | |
| Grade 11 | 445 | |
| Grade 12 | 596 | |

## Course Sale Infomation

| | Avg Price | Course Qty | Sold Qty | Pct of Sold |
|---|---|---|---|---|
| Biology | 953 | 2,830 | 1,493 | 23.6% |
| Chemistry | 1,003 | 6,850 | 4,135 | 26.4% |
| English | 983 | 30,268 | 13,669 | 24.4% |
| Geography | 999 | 2,230 | 743 | 14.8% |
| History | 977 | 2,290 | 1,102 | 20.4% |
| Literature | 936 | 32,529 | 15,694 | 25.1% |
| Math | 980 | 77,836 | 52,104 | 40.2% |
| Physics | 1,004 | 11,870 | 5,463 | 21.9% |
| Politics | 971 | 1,810 | 455 | 10.1% |

## Course Revenue By Grade



- 5.76%
- 8.90% Grade 4
- 7.68% Grade 5
- 21.52% Grade 10
- 9.60% Grade 1
- 9.36% Grade 2
- 7.63% Grade 6
- 6.77% Grade 9
- 10.61% Grade 3
- 8.01% Grade 7
- 5.12% Grade 8

# China's K12 Online Education Data Management System

## Top 3 course each grade ranked by revenue

| Grade Id | Course Id | Course Name | | Revenue |
|---|---|---|---|---|
| 1 | 3153689 | 一年级语文春季系统班(周.. | | 108,000 |
| | 3168085 | 【目标A班】一年级数学春季.. | | 108,000 |
| | 3157973 | 【目标A+班】一年级剑桥英语.. | | 117,450 |
| 2 | 3153833 | 二年级语文春季系统班(周.. | | 108,000 |
| | 3153837 | 二年级语文春季系统班(周.. | | 108,000 |
| | 3168069 | 【目标A班】二年级数学春季.. | | 108,000 |
| 3 | 3158007 | 【目标M班】三年级剑桥英语.. | | 113,100 |
| | 3158011 | 【目标M班】三年级剑桥英语.. | | 117,450 |
| | 3158039 | 【目标A班】三年级剑桥英语.. | | 121,800 |
| 4 | 3153739 | 四年级语文春季系统班(周.. | | 108,000 |
| | 3169421 | 【目标A班】四年级数学春季.. | | 108,000 |
| | 3169425 | 【目标A班】四年级数学春季.. | | 108,000 |
| 5 | 3153747 | 【目标A班】五年级数学春季.. | | 104,400 |
| | 3153817 | 五年级语文春季系统班(周.. | | 108,000 |

## Top 3 course each subject ranked by revenue

| S.. | Subject N.. | Course Id | Course Name | | Revenue |
|---|---|---|---|---|---|
| 1 | Politics | 3893787 | 新高二政治秋季系统班(周.. | | 151,200 |
| | | 3893789 | 新高二政治秋季系统班(周.. | | 142,800 |
| | | 3881337 | 新高二政治暑期系统班(暑3.. | | 64,800 |
| 2 | Geography | 3895105 | 新高二地理秋季系统班(周.. | | 155,400 |
| | | 3162269 | 高一地理春季系统班(周日0.. | | 125,874 |
| | | 3162267 | 高一地理春季系统班(周六1.. | | 122,877 |
| 3 | History | 3905257 | 新高二历史秋季系统班(周.. | | 147,000 |
| | | 3163151 | 高一历史春季系统班(周六2.. | | 131,868 |
| | | 3163149 | 高一历史春季系统班(周六1.. | | 131,868 |
| 4 | Biology | 3878607 | 【目标985班】新高二生物秋.. | | 151,200 |
| | | 3164481 | 【目标一本班】高一生物春季.. | | 146,853 |
| | | 3878603 | 【目标985班】新高二生物秋.. | | 142,800 |
| 5 | Chemistry | 3876275 | 【目标985班】新高二化学.. | | 184,800 |
| | | 3876407 | 【清北高端班】新高二化学秋.. | | 163,200 |

## Top10 course by longest sale time(days)

| Index | Course Id | Sale Length |
|---|---|---|
| 1 | 3877605 | 163 |
| 2 | 3878911 | 163 |
| 3 | 3878913 | 163 |
| 4 | 3891395 | 163 |
| 5 | 3927337 | 163 |
| 6 | 3927349 | 163 |
| 7 | 3927517 | 163 |
| 8 | 3875255 | 163 |
| 9 | 3875603 | 163 |
| 10 | 3875605 | 163 |

Sale Length