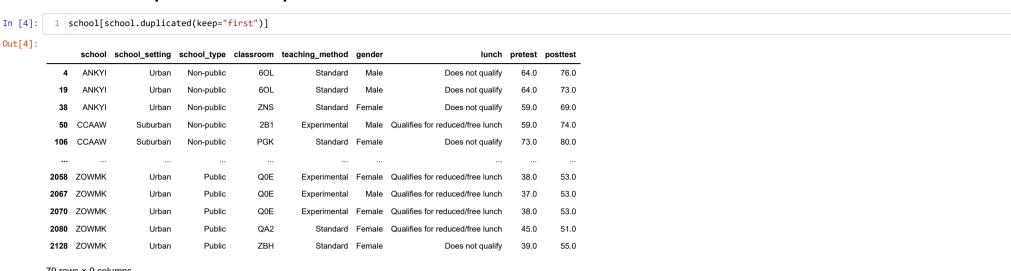
To import necessary modules

```
In [1]: 1 import pandas as pd import numpy as np import warnings warnings.filterwarnings("ignore")

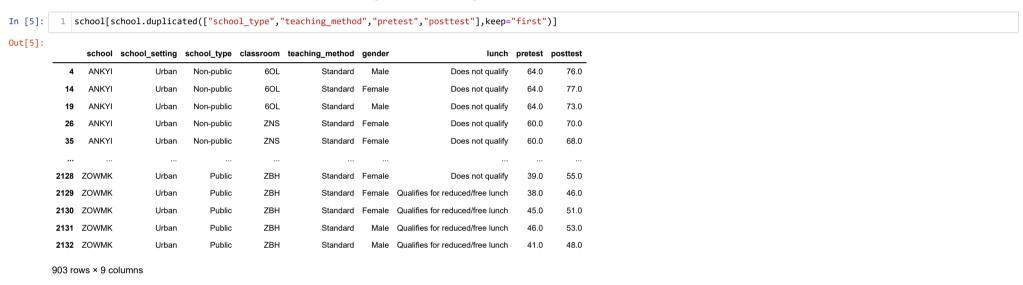
In [2]: 1 test=pd.read_csv(r"K:\Desktop\NIIT\tables\DS1_C4_S4_Test_Scores_Data_Practice.csv") non_unique=test.columns[test.columns!="student_id"]

In [3]: 1 school=test[non_unique]
```

Task1: To help find all the duplicate records after thier first occurence



Task2: To find all duplicates of school type teaching method pre test and post test parameters



Task3: To find average pretest score for all schools and highest pretest score

```
1 | avg=school.groupby(["school"]).mean()["pretest"]
In [6]:
          2 print(avg)
          4 print("The school with the highest average pretest is = ",avg[avg==avg.max()])
        school
        CCAAW
        CIMBB
                 65.067568
                 53.925234
        CUQAM
                 54.327869
        DNQDD
        FBUMG
                 62.891304
        GJJHK
                 53.194915
        GOKXL
                 50.796875
        GOOBU
                 38.248408
        IDGFP
                 75.202128
                 41.865385
        KFZMY
        KZKKE
                 37.261261
        LAYPA
                 62.035088
                 52.527027
        UAGPU
                 62.390805
        UKPGS
                 78.453125
        UUUQX
                 67.253012
        VHDH
                 52.666667
        VVTVA
                 35.955752
                 68.130435
        ZMNYA
        ZOWMK
                41.572650
        Name: pretest, dtype: float64
        The school with the highest average pretest is = school
                78.453125
        Name: pretest, dtype: float64
```

Task4: To display the location wise and gender wise maximum

Task5: To display the minimum and maximum post test scores for all the schools and also show the top five post test scores Top 5

```
posttest_sc=pd.pivot_table(school,index="school",values=["posttest"],aggfunc=["min","max"])
In [8]:
             print("The minimum and maximum posttest scores by all schools")
        The minimum and maximum posttest scores by all schools
Out[8]:
                 posttest posttest
          school
          ANKYI
                    63.0
                            79.0
         CCAAW
                    67.0
                            91.0
          CIMBB
                    64.0
                            88.0
                    56.0
                            76.0
         DNQDD
                    49.0
                            79.0
         FBUMG
                    68.0
                            88.0
          GJJHK
                    49.0
                            85.0
         GOKXL
                    48.0
                            77.0
         GOOBU
                    32.0
                            64.0
          IDGFP
                    74.0
                            100.0
          KFZMY
                    44.0
                            67.0
          KZKKE
                    36.0
                            62.0
          LAYPA
                    63.0
         OJOBU
                    50.0
                            84.0
         QOQTS
                    51.0
                            85.0
                    62.0
                            83.0
         UKPGS
                    82.0
                            99.0
         UUUQX
                    62.0
                            91.0
          VHDHF
                    52.0
                            77.0
         VKWQH
                    48.0
                            82.0
          VVTVA
                    39.0
                            62.0
          ZMNYA
                    66.0
                            95.0
         ZOWMK
                    43.0
                            63.0
In [9]:
            print("The top 5 posttest scores ")
             m=pd.pivot_table(school,index="school",values=["posttest"],aggfunc=["max"]).sort_values(by=[('max', 'posttest')],ascending=False).iloc[0:5,:]
        The top 5 posttest scores
Out[9]:
                 posttest
          school
          IDGFP
                   100.0
         UKPGS
                    99.0
         ZMNYA
                    95.0
         UUUQX
                    91.0
         CCAAW
                    91.0
```

Task6: To display school wise averages and sum of pre test and top 5 pre test scores

ZOWMK

41.572650 52.905983

```
print("The school wise average of pretest and postest")
print(school.groupby("school").mean())
In [10]:
         The school wise average of pretest and postest
                   pretest posttest
          school
         ANKYI
                 61.341463 71.390244
         CCAAW
                 64.623853 78.110092
         CIMBB
                 65.067568 76.945946
                 53.925234 65.560748
         CUQAM
         DNQDD
                 54.327869 66.565574
         FBUMG
                 62.891304 78.608696
                 53.194915 65.025424
         GOKXL
                 50.796875 64.953125
                 38.248408 49.643312
          GOOBU
         IDGFP
                 75.202128 87.223404
         KFZMY
                 41.865385 54.576923
                 37.261261 47.918919
          KZKKE
          LAYPA
                 62.035088 73.571429
                 56.197531 67.814815
         OJOBU
         QOQTS
                 52.527027 64.671141
         UAGPU
                 62.390805 71.873563
         UKPGS
                 78.453125 91.188976
                 67.253012 79.261905
         UUUQX
         VKWQH
                 52.060000 64.820000
                 35.955752 49.203540
         VVTVA
         ZMNYA
                 68.130435 81.608696
```

```
In [11]:
             print("The sum of pretest scores : ")
             print(school.groupby("school").sum()["pretest"])
             print("The top 5 average pretest scores = ")
             print(school.groupby("school")["pretest"].mean().sort_values(ascending=False).head())
         school
                   2515.0
         ANKYI
         CCAAW
                   7044.0
         CIMBB
                   4815.0
         CUQAM
                   5770.0
         DNQDD
         FBUMG
                   2893.0
         GJJHK
                   6277.0
         GOKXL
                   3251.0
         GOOBU
                   6005.0
         IDGFP
                   7069.0
         KFZMY
         KZKKE
                   4136.0
         LAYPA
                   3536.0
         OJOBU
                   4552.0
                   7774.0
         QOQTS
         UAGPU
                   5428.0
         UKPGS
         UUUQX
         VHDHF
                   2686.0
         VKWQH
                   5206.0
                   4063.0
         VVTVA
         ZMNYA
                   4701.0
                   4864.0
         Name: pretest, dtype: float64
         The top 5 average pretest scores =
         school
                  78.453125
         UKPGS
         IDGFP
                  75.202128
         ZMNYA
                  68.130435
         UUUQX
                  67.253012
         CIMBB
                  65.067568
         Name: pretest, dtype: float64
```

Task7: The school FBUMG with type of 5LQ classroom wishes to offer scholarships

```
In [12]:
               FBU=test[(test.school=="FBUMG") & (test.classroom=="5LQ")].loc[:,["student_id","posttest"]]
               FBU["Scholarship"]=(FBU.posttest/100)*10000
Out[12]:
               student_id posttest Scholarship
                  04DG5
                                      8400.0
           467
                             84.0
           468
                  20M2D
                             85.0
                                       8500.0
           469
                  39KCW
                             81.0
                                      8100.0
           470
                   5Z1B6
                             88.0
                                      8800.0
           471
                   6TLU8
                             81.0
                                      8100.0
           472
                  ATQQJ
                             82.0
                                      8200.0
           473
                  AYEU1
                             79.0
                                      7900.0
           474
                   B9FSU
                             83.0
                                       8300.0
           475
                   EV13K
                             83.0
                                      8300.0
           476
                   15H37
                             83.0
                                      8300.0
           477
                    JC519
                             80.0
                                       0.0008
           478
                             81.0
                                      8100.0
                   JPE2J
           479
                  MFBYU
                             82.0
                                      8200.0
           480
                   O144X
                             86.0
                                       8600.0
           481
                  OGKP3
                             83.0
                                      8300.0
           482
                   UMFI7
                             73.0
                                      7300.0
           483
                   V1DNJ
                             76.0
                                       7600.0
           484
                  YRN9S
                                      8700.0
```

Task8: To replace does not qualify and qualifies with reduced free lunch with NO and YES

```
In [13]:
             school.lunch.replace(to_replace=["Does not qualify","Qualifies for reduced/free lunch"],value=["No","Yes"],inplace=True)
             school.lunch
Out[13]: 0
                  No
                  No
                  No
                 No
                 Yes
         2130
         2131
                 Yes
         2132
                 Yes
         Name: lunch, Length: 2133, dtype: object
```

Task9: To identify count of teaching methods by locations

```
pd.crosstab(school.school_setting,school.teaching_method)
In [14]:
Out[14]:
          teaching_method Experimental Standard
            school_setting
                                        309
                                284
                Suburban
                                        433
                   Urban
                                275
                                         631
              school.value_counts().groupby(["school_setting","teaching_method"]).sum()
In [15]:
Out[15]: school_setting teaching_method
                         Experimental
                                             199
         Rural
                         Standard
                                             307
         Suburban
                          Experimental
                                             283
         Urban
                          Experimental
                                             269
                          Standard
                                             626
         dtype: int64
```

Task10: To replace the pre test scores of male students with new value of 135 for male students who have post test less than 39

```
index=school[(school.gender=="Male")&(school.posttest<39)]["pretest"].index
school.pretest[index.values]=135
school.pretest[index.values]</pre>
In [16]:
                school[(school.gender=="Male")&(school.posttest<39)]</pre>
Out[16]:
                   school_setting school_type classroom teaching_method gender lunch pretest posttest
             742 GOOBU
                                                            HKF
                                                                                                               38.0
                                   Urban
                                                Public
                                                                          Standard Male
                                                                                              Yes
                                                                                                   135.0
             745 GOOBU
                                   Urban
                                                Public
                                                            HKF
                                                                          Standard Male
                                                                                              Yes
                                                                                                     135.0
                                                                                                               38.0
```

Task11: To compare average posttest scores with sub urban and urban male students

Yes

Yes

135.0

135.0

135.0

135.0

34.0

37.0

38.0

37.0

Standard Male

Standard Male Yes

Standard Male

Standard Male Yes

Task12: To help replace missing values with most recurring values

Public

Public

Public 3D0

Public QTU

Urban

Rural

Rural

Rural

751 GOOBU

1023 KZKKE

1081 KZKKE

1089 KZKKE

HKF

```
In [18]:
             print(school.isnull().sum())
             school.gender=school.gender.fillna("Male")
             school.pretest=school.pretest.fillna(school.pretest.mean())
             school.posttest=school.posttest.fillna(school.posttest.mean())
             print("After replacing ")
             print(school.isnull().sum())
         school\_setting
         school_type
                             0
                             0
         classroom
         teaching_method
                             0
         gender
                            10
         pretest
         posttest
         dtype: int64
         After replacing
         school\_setting
         school_type
         classroom
         {\tt teaching\_method}
         gender
         pretest
         posttest
         dtype: int64
```