# Data Warehousing, Integration and Provenance-CS520

## Vizier Assignment
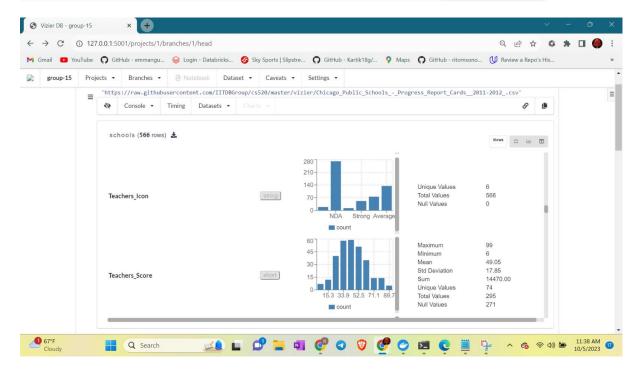
## Group 15

Rohith Reddy Bairi          rbairi@hawk.iit.edu

Deekshitha Tummala          dtummala@hawk.iit.edu

Shirisha Vaddegoni          svaddegoni@hawk.iit.edu
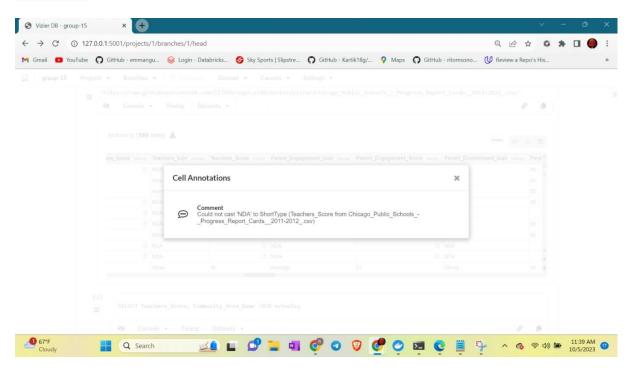
Brijesh Kumar Yadav Bandi   bbandi@hawk.iit.edu

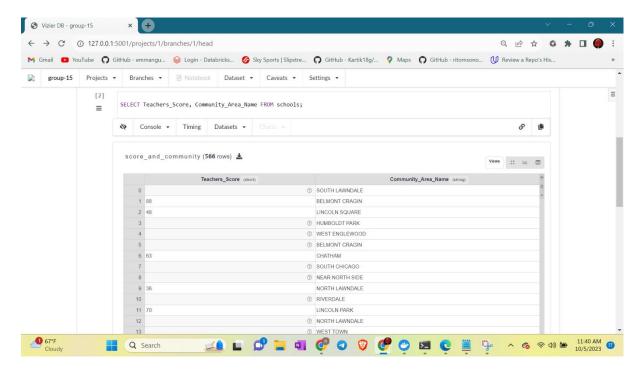**Task-1:** Load a dataset and take a screenshot of the result.

**Task-2:** Select the detail view and look at the distributions of some columns. Then look at the column view and take a screenshot of the distribution for column Teachers_Score.
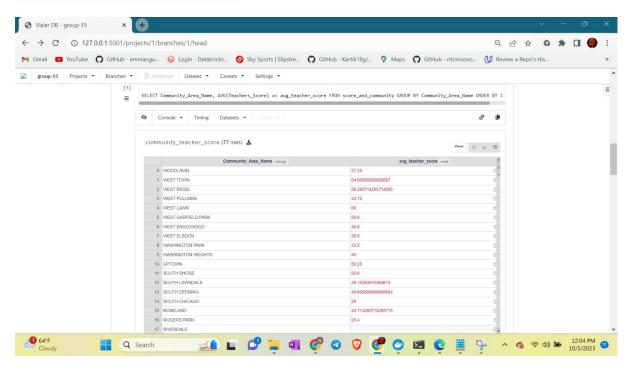


**Task-3:** Click on one of the question marks for values in the teachers column and take a screenshot.

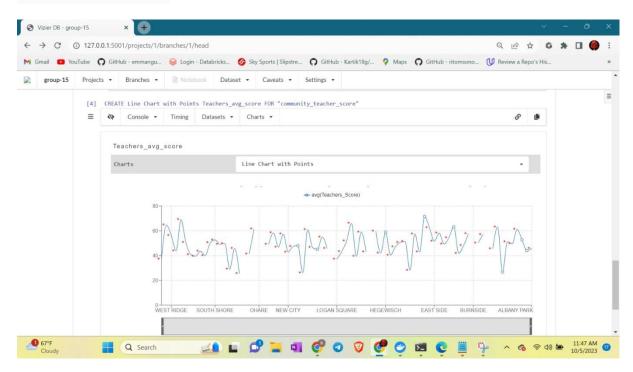**Task-4:** Create a SQL cell and write a query that returns columns Teachers_Score and Community_Area_Name. SQL results can be stored as new dataset score_and_community. And take a screenshot of the result.
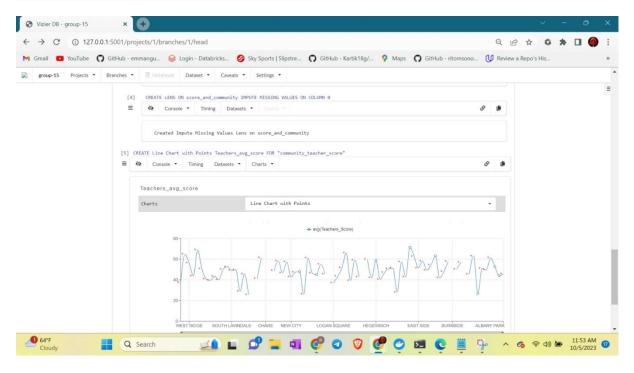


**Task-5:** Create a SQL cell and write a query over the score_and_community dataset that computes the result as described above. Call the result dataset community_teacher_scores. And take a screenshot of the result.

**Task-6:** Create a line chart of the aggregation result by creating a plot cell and take a screenshot of the result.



**Task-7:** Insert a new cell above the SQL cell that computes the average teacher scores (notebooks in Vizier are executed top down) by pressing the three bars below the cell number. Select *"Impute Missing Values"*, select the score_and_community dataset and Teachers_Score as the column to be imputed, and select mean as the imputation method and take a screenshot of the updated line chart.

**Task-8:** Create a Python cell at the end of the notebook and create a function called print_avg_teachers that uses Vizier's API to get a handle for this dataset and print all values of the avg_teacher_score column. Hint: Use the "Show Code Examples" button to see example Vizier API usage and see here for the API documentation. Then use vizierdb.export_module to export the function. Then create a second Python cell and use vizierdb.get_model("print_avg_teachers") for importing the function and then call it. Take a screenshot of the result.



**Task-9:** Create another Python cell and use Vizier's API to access the dataset community_teacher_scores as a DataFrame, then filter out rows where the avg_teacher_score is larger than or equal to 30.0 and then print the remaining rows and take a screenshot.

group-15   Projects ▾   Branches ▾   Notebook   Dataset ▾   Caveats ▾   Settings ▾

```
61.25
52.66666666666664
44
45.75
```

[8]

```python
df = vizierdb.get_data_frame('community_teacher_score')
print(df.loc[df['avg_teacher_score'] > 30.0])
```

Console ▾    Timing    Datasets ▾    Charts ▾

```
   Community_Area_Name  avg_teacher_score
0             WOODLAWN          37.250000
1            WEST TOWN          64.666667
2           WEST RIDGE          56.285714
3         WEST PULLMAN          43.750000
4            WEST LAWN          69.000000
..                 ...                ...
72      AUBURN GRESHAM          49.571429
73             ASHBURN          61.250000
74       ARMOUR SQUARE          52.666667
75      ARCHER HEIGHTS          44.000000
76         ALBANY PARK          45.750000

[65 rows x 2 columns]
```

+