# Fall 2022 - CS 673:  Scalable Databases

## Midterm Project (Team Submission) Requirements

Total Points : 100

Due Date: 11/04/2022

| Task | Points | Self-assessment |
|---|---|---|
| **Problem Definition:**<br>What is the problem you are working on?<br>Which dataset are you using? (Select structured and unstructured ) | 10 | |
| **PowerPoint Presentation :**<br>Your mid-term presentation must have Introduction, Problem definition or hypothesis, dataset used, your EDA backed by appropriate visualizations, results, and conclusion. | 30 | |
| **Python Notebook:**<br>Use of python / SQL concepts taught in the class such as functions, creating own packages, file Handling, exception Handling, use of third-party libraries. Demonstrate your knowledge on handling structured and unstructured data.<br><br>**Essential key steps to demonstrate in your Python Notebook**<br>1. Loading data in to DataFrames (both structured and unstructured). Integration of SQL and Python<br>2. Check the Data Types of your data columns.<br>3. Drop any NULL, missing values or unwanted columns.<br>4. Drop duplicate values.<br>5. Check for outliers using a box plot or histogram.<br>6. Plot features against each other using a pair plot.<br>7. Use a HeatMap for finding the correlation between the features(Feature to Feature).<br>8. Use a scatter plot to show the relationship between 2 variables.<br>9. Merging two Data Frames.<br>10. Slicing Data of a particular column value (like year, month, filter values depending on the categorical data)<br>11. Representing data in matrix form.<br>12. Upload data to Numerical Python (NumPy)<br>13. Select a slice or part of the data and display.<br>14. Use conditions and segregate the data based on the condition (like show data of a feature(column) >,<,= a number) | 50 | |

| | | |
|---|---|---|
| 15. if applicable Use mathematical and statistical functions using libraries.<br>16. Select data based on a category(categorical data based).<br>17. Libraries expected to try(minimum 4 required): Pandas, Numpy, Seaborn, Matplotlib .<br>18. Write your own functions and handle exceptions in the functions.<br>19. Use of *arg and **kwargs.<br>20. Use of data functions. | | |
| **Teamwork:**<br> Effective communication and participation with your teammates. | 10 | |

## Submission :

- Submit the all the files with code, data set, and a word document or PPT with the explanation individually.
- You are required to do a class presentation, with all team members participating. The presentation will be on zoom, you have switch on the camera.
- Late submission up to one week, 20% deduction of total points earned.
- Submit the self-assessment along with the above-mentioned files.

*Important : No plagiarism, please implement your own idea and submit your own work. Your work will be checked for plagiarism.*

*NOTE : Extra*

*credit for the class participation (Q&A)  – 5 points.*