# CSCI S-101 Python for Engineering
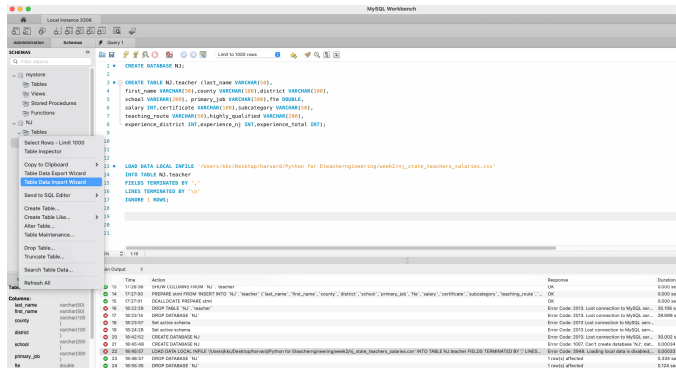
<div align="right">Name: Qi Chu</div>

## Data Uploading in SQL and Python

For MySQL uploading, I used mysql.connector in Jupyter notebook. It can also be done in MySQL WorkBench through Table Data Import Wizard:

 Or in terminal, using the same codes.

Python uploading was just using pd.read_csv.

## Data Cleaning

1. I first check the data types, and observe most of the variables are identified as objects, so will need to assign the new, appropriate data types to them.
2. Rows with any NaN value for any variable are dropped, even some data entries have only one NaN value, it is not possible to use the group average or previous/next available value since every teacher's information is independent. The NaN values do not add value to the data analysis.
3. The variables that contain text values, e.g. last_name, primary_job, are converted to strings, and are trimmed with leading and trailing spaces.
4. For the variables that are supposed to be numeric, I first check if there's any non-numeric values hiding by forcing them into float values, any errors will be replaced with NaN values.
5. Step 5 identified another 10 non-numeric / NaN values, so I drop the rows with NAs again to have the final clean data. Converting the numeric variables into float is only to be safer, but I also check if separating by float and integer is necessary. I re-convert these float variables into strings, and strip out the last character, if it's 0, that original value is an integer, otherwise it's a real float. By doing so, I figure out only 'fte' is a float while the rest variables are integers. This validates the SQL inputs too.
6. Biggest Challenge:  Python identifies most variables as objects, then data reading, comparing, trimming are not possible, will need to convert to appropriate data types. There are also hidden non-numeric values that need to be cleaned.