# Data Warehouse & Data Mining Lab

Submitted by

```
1  Gyanendra Shukla
2  CSE 1
3  191112040
```

## Assignment Problem

Given the CSV file with field descriptions, convert the data to a structured MySQL table.

## Approach Used

I've used SQL to create and import data in an MySQL table. In MySQL there is a load data method that allows us to load some data with filtering on it. First of all, I created a new table with the given schema and observing the data. For field separation, I separated them with `,` and for new entry separation with `'\n'` . I have finally displayed first 10 entries of the dataset.

## Code

```sql
1
2  -- Creating a MySQL table according to the given schema.
3  -- The fields that were continuos were marked as float, and other
4  -- fields with fixed values as VARCHARs.
5
6  CREATE TABLE socialinfo (
7      id INT NOT NULL auto_increment,
8      age float,
9      workclass VARCHAR(100),
10     fnlwt float,
11     education varchar(100),
12     educationnum float,
13     maritalstatus varchar(100),
14     occupation varchar(100),
15     relationship varchar(100),
16     race varchar(100),
17     sex varchar(100),
18     capitalgain float,
19     capitalloss float,
20     hoursperweek float,
21     nativecountry varchar(100),
22     salary varchar(100),
23     primary key (id)
24 );
25
26
```

```
27    -- After creating the table, we load the local infile and store it in our
      table
28    load data local infile 'C:/ProgramData/MySQL/MySQL Server
      8.0/Uploads/Sample.txt'
29    into table socialinfo
30    fields terminated by ','
31    lines terminated by '\n'
32    (age, workclass, fnlwt, education, educationnum, maritalstatus, occupation,
      relationship,
33    race, sex, capitalgain, capitalloss, hoursperweek, salary);
34
35    -- showing the top 10 results
36    select * from socialinfo limit 10;
37
38    -- showing number of entries in the dataset
39    select count(*) from socialinfo;
```

## Description of Code

We can create a new MySQL table using CREATE TABLE method. I've created the `socialinfo` table with appropriate schema. I've also added an extra `id` field that is auto incrementing and set that as the primary key. Then, I loaded the `Sample.csv` file and inserted it in the table. For separating the data we're using `,` for fields and `\n` for entries.
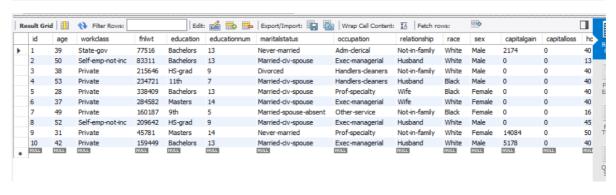
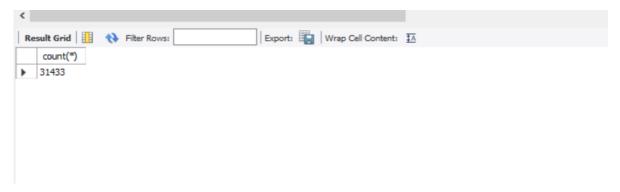## Output Snapshots



Fig: Top 10 entries of the dataset



Fig: Total entries in dataset

1. assignment problem, 2. approach used, 3. code with proper comments, 4. description of code, 5. output snapshots.