

# Phase 2: Database

Sierra Allred

July 25, 2016

The second phase of my thesis was the creation of a MySQL database to organize data collected from College Scorecard and AreaVibes.com. Because the visualization portion of my project was done with a team before I had database experience, the visualization and database are not connected at this time. The next step of my project would be to use Node.js to connect the visualization to the database with the queries I have created.

This document will describe how I collected the data I needed from AreaVibes, how the project database is organized, and the queries that I created to collect data for the visualization and statistical analysis portions of my project.

## 1 Collection of AreaVibes Data

For the statistical analysis portion of my project I needed cost of living (COL) data for various cities throughout the United States. AreaVibes.com gives each city a cost of living rating, but their data was not readily available from an API, and was spread throughout the site on various web pages. I used Scrapy, a Python based web crawling language, to collect the data I needed in an organized csv file.

I started by creating a Scrapy item that defined all the fields I wanted to collect. The COLarea was the COL of a certain area of a city, and COLcity was the overall COL of the city. I also collected the City, Area, and State.

```
from scrapy.item import Item, Field

class ColItem(Item):
    # define the fields for item:
    COLarea= Field()
    COLcity= Field()
    City = Field()
    State = Field()
    Area = Field()
```

Next I created the spider I would use to crawl the website. I went to AreaVibes robots.txt file to find the sitemap, and use that to navigate to a xml file containing all the links that held data on each city and its COL. I used that xml as the starting url for my spider to crawl. My spider collects the page from each url, and then creates an item to store the data it collects. I used html tags to direct the spider to the data in the page, which it then extracts and stores.

```

import scrapy
from scrapy.spiders import SitemapSpider
from COL.items import ColItem

class COLSpider(SitemapSpider):
    name = "COLspidey"
    sitemap_urls = ['http://www.areavibes.com/sitemap/sm_hood_COL.xml']

    def parse(self, response):
        page = response.xpath("//html")
        for sel in page:
            item = ColItem()
            item["title"] = sel.xpath("//title/text()").extract()
            item["COLarea"] = sel.xpath("//tr[@class='summary major']/td[2]/text()").extract()
            item["COLcity"] = sel.xpath("//tr[@class='summary major']/td[3]/text()").extract()
            item["City"] = sel.xpath("//nav[@class='breadcrumbs']/a[3]/text()").extract()
            item["State"] = sel.xpath("//nav[@class='breadcrumbs']/a[2]/text()").extract()
            item["Area"] = sel.xpath("//nav[@class='breadcrumbs']/a[4]/text()").extract()

            yield item

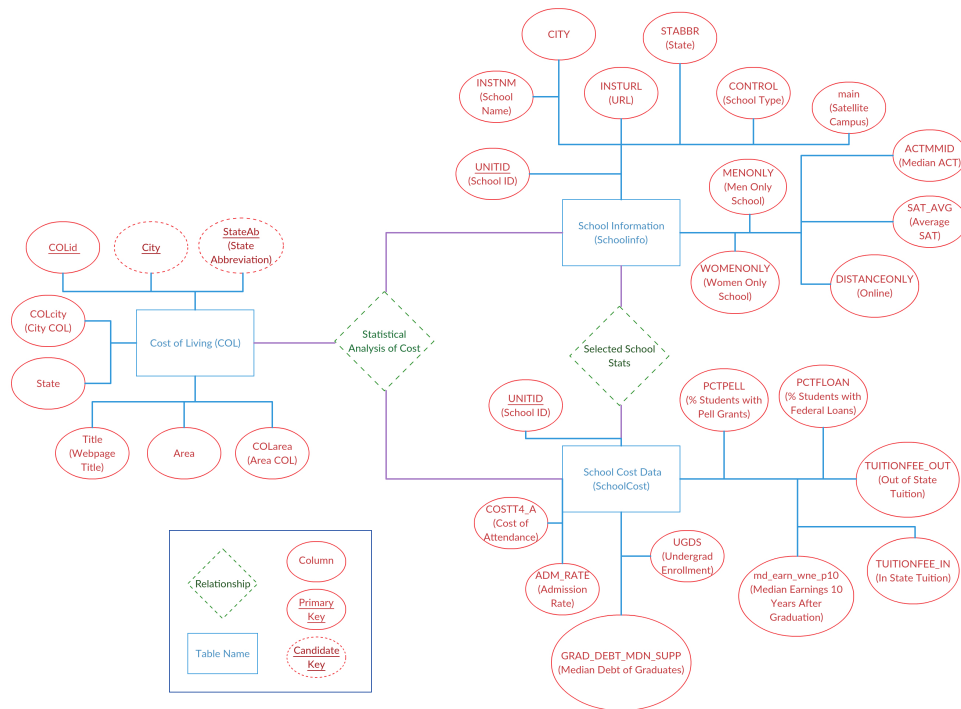
```

I used the command line to crawl my spider and store all the data in an organized csv file which could easily be added to the project database.

## 2 Database Organization

The project database is a very simple MySQL database made up of 3 tables. One table stores all of the cost related data about each university in the United States, and another holds all other relevant information about each school such as the location, url, name, and test scores. The last table holds all of the COL data collected from AreaVibes.

Below is the entity-relationship model for the project database outlining the tables, columns, and keys and their relationships.



### 3 Queries

I have created all the queries needed for the visualization and statistical analysis portions of my project, this section will describe each of these.

The create table statement used for the COL table.

```
1 CREATE TABLE 'scorecard'.'COL' (  
2   'COLid' INT NOT NULL,  
3   'City' VARCHAR(45) NULL,  
4   'COLcity' INT NULL,  
5   'COLarea' INT NULL,  
6   'StateAb' VARCHAR(45) NULL,  
7   'Title' VARCHAR(45) NULL,  
8   'Area' VARCHAR(45) NULL,  
9   'State' VARCHAR(45) NULL,  
10  PRIMARY KEY ('idCOL'),  
11  UNIQUE('City', 'State'));
```

This statement was used to load the COL data in the csv into the COL table.

```
1 load data local infile '/home/sierra/Desktop/COLdata.csv'  
2 into table COL fields terminated by ','  
3 enclosed by '"'  
4 lines terminated by '\n'  
5 IGNORE 1 LINES  
6 (idCOL, City, COLcity, COLarea, StateAb, title, Area, State);
```

This query was used to collect the cost data I needed for statistical analysis. I had to join all 3 tables by School ID, City, and State. Because so much data was missing from the College Scorecard dataset, I specified that the query only return rows that did not have NULL or PrivacySuppressed values.

```
1 select DISTINCT  
2   cols.COLcity COL,  
3   sc.COSTT4_A cost,  
4   sc.md_earn_wne_p10 earnings,  
5   sc.GRAD_DEBT_MDN_SUPP debt  
6 from  
7   SchoolCost sc,  
8   COLs cols,  
9   Schoolinfo si  
10 where  
11   si.STABBR = cols.StateAb  
12   and si.CITY = cols.City  
13   and si.UNITID = sc.UNITID  
14   and sc.COSTT4_A not like 'NULL'  
15   and sc.md_earn_wne_p10 not like 'NULL'  
16   and sc.md_earn_wne_p10 not like 'Priva'  
17   and sc.GRAD_DEBT_MDN_SUPP not like 'NULL'  
18   and sc.GRAD_DEBT_MDN_SUPP not like 'PrivacySuppressed';
```

These queries create temporary tables. schoolList holds the name and id of all the schools returned by the user driven filter. All the variables contained by [] are placeholders for values returned by the filter. graph selects all the data needed for the different options of the bar chart in the visualization from schoolList.

```

1  where
2  schoolList as (
3  select
4      si.UNITID sid,
5      si.INSTNM sname
6  from
7      SchoolCost sc,
8      SchoolInfo si
9  where
10     sc.UNITID = si.UNITID
11     and si.SAT_AVG < [satHi]
12     and si.SAT_AVG > [satLo]
13     and si.ACTCMMID < [actHi]
14     and si.ACTCMMID > [actLo]
15     and sc.TUITIONFEE_OUT < [tutHi]
16     and sc.TUITIONFEE_OUT > [tutLo]
17     and sc.TUITIONFEE_IN < [tutHi]
18     and sc.TUITIONFEE_IN > [tutLo]
19     and sc.UGDS < [sizeHi]
20     and sc.UGDS > [sizeLo]
21     and si.CONTROL = [type]
22 ),
23 graph as (
24 select
25     dat.ADM_RATE admin,
26     dat.COSTT4_A cost,
27     dat.md_earn_wne_p10 earnings,
28     dat.GRAD_DEBT_MDN_SUPP debt
29 from
30     SchoolCost sc,
31     SchoolList sl
32 where
33     sl.UNITID = sc.UNITID
34 )

```

This query is for the part of the visualization that gives you information about an individual school selected from the list of filtered schools. [selectedSchool] is a placeholder for ID of the school selected from that list.

```

1  select
2      si.INSTNM name,
3      si.CITY city,
4      si.STABBR state,
5      si.INSTURL url,
6      si.MENONLY men,
7      si.WOMENONLY women,

```

```

8      si.main satellite,
9      si.DISTANCEONLY onlineOnly,
10     sc.ADM_RATE admin,
11     sc.UGDS undergradEnroll,
12     si.ACTCMMID medACT,
13     si.SAT_AVG avgACT,
14     sc.TUITIONFEE_IN inTuition,
15     sc.TUITIONFEE_OUT outTuition,
16     sc.COSTT4_A cost,
17     sc.PCTPELL pell,
18     sc.PCTFLOAN loan,
19     sc.GRAD_DEBT_MDN_SUPP debt,
20     sc.md_earn_wne_p10 earnings
21 from
22     SchoolCost sc,
23     SchoolInfo si
24 where
25     si.UNITID = [selectSchool]
26     and si.UNITID = sc.UNITID;

```

This query provides the data needed for the selected schools list to be ordered by. [selectOrder] is a placeholder for the option selected to order the list.

```

1  select
2      si.UNITID sid,
3      si.INSTNM name,
4      sc.UGDS undergradEnroll,
5      sc.GRAD_DEBT_MDN_SUPP debt,
6      sc.COSTT4_A cost,
7      sc.ADM_RATE admin,
8      sc.md_earn_wne_p10 earnings,
9  from
10     SchoolCost sc,
11     SchoolInfo si,
12     SchoolList sl
13 where
14     sc.UNITID = sl.UNITID
15     and sc.UNITID = si.UNITID
16 order by
17     sc.[selectOrder];

```