



# **Normalized Football Database**

**IS426**

**By Jackson Nahom**

## **Introduction:**

For this project, I took it upon myself to make a normalized database of historical soccer data that can have its data analyzed. To do this I had to find a data source to get the historical data. This source came to be realized at football-Data.co.uk. This website is mostly used for analyzing soccer for the purpose of betting, it also has the historical betting odds. However, the betting odds are not what I am looking to analyze. On Football Data there are CSV s for leagues in 27 countries.

## **Building The Database:**

To start off I looked at the CSVs in England link on football-Data.co.uk, and realized it would be simplest to web scrape all the CSVs automatically. To do this I had the unique sequence on the lines with the CSV links. For England, this turned out to be: <IMG SRC="Excel. From there I isolated the CSV link using the split function. From there I can use a base URL and the link I split out to use request.get to take the data in the CSV. From there I can write it to a file that I name in this format YEAR\_CD.csv, where C is country and D is division. Now that I could download one country's CSV files, the next step was to go through all the countries. This was a similar problem as the last step. I had to parse through and look for the country URL and add it to a base URL to use requests.get to the parse through that webpage using the code I tested on the England page. I hit a roadblock here because Football Data split the leagues into two formats, main leagues, and extra leagues. The way the CSV files are stored in the extra leagues than in main leagues. To get around this I found a sequence that was unique to the lines with CSV files in extra leagues that did not appear in the source page for the

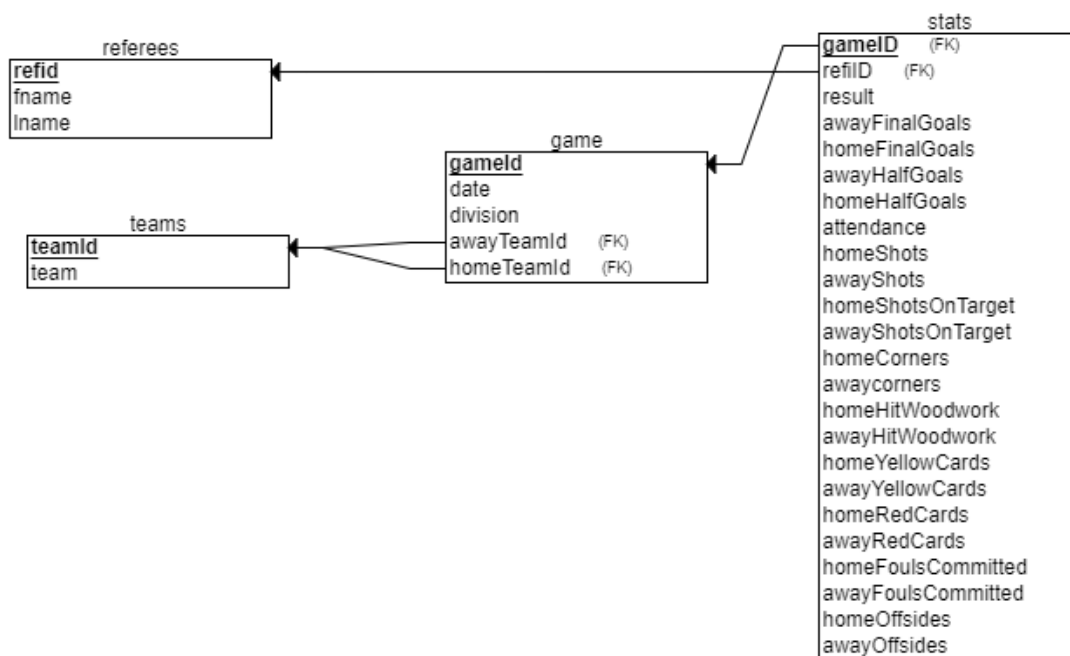
main leagues. This turned out to be: `.csv"><IMG`. One other thing I did to prevent errors was to replace `/` with `_` when saving files because you can't name files with slashes. When this script was finished all in all it downloaded 551 CSV files from 27 countries. For the code for this section see `downloadCSV.py` in Appendix.

Now That I had collected all the CSVs I had to upload it to a source database. At first, I had errors trying to upload CSV files one at a time. This was because the columns did not match every CSV. So I wrote a script to go through all the CSV files that I wanted to upload to go through getting all the headers and get rid of duplicates and make the source table with nothing in it. I used the `glob` dictionary to call `glob.glob` to open up all the CSV files. I needed to populate the table with a different script. To see the create table script see `createSourceTable.py` in Appendix.

Now that I had created the table with all the headers needed in CSVs that were to be uploaded. For the analysis's sake of the project, I only uploaded the data from the first division of England. I wrote a function called `upload` that goes through a csv and puts the data in the source file. This function is very similar to the one we did in class with a few changes. I used the `glob.glob` function again to collect all the CSV files. The data from England's first division went back to 1993 and resulted in 9804 entries to the source table. The two scripts necessary to achieve this is `Upload Function` and `uploadSource.py` in the Appendix. A couple problems that I ran into while going through this process were that a couple CSV files returning blank headers or more rows than in the CSV. I got around these problems by checking if there is anything in the header. I also only uploaded a row if the date existed because there was always a date. After

solving that I had to fix the date format so that it was in the SQL format of YYYY-MM-DD. The way that Football Data had date was DD/MM/YYYY. To fix this I wrote a function to convert the date. This is located in the Appendix as formatDate Function.

Once I have the source table I can Normalize the data. The first step of that is creating a normalized relational schema. This led me to make this relationship schema:



Now that I have my schema I created these table. For the referee's table, none of the values can be NULL. The same goes for the teams, and game tables. However, for the stats table the only values that will never be NULL are gamelID, result, awayFinalGoals, and homeFinalGoals. The rest of the values in stats can be NULL if necessary. To populate the referees, and teams table I had to put all the names for referees and teams into a list, then I could go through and get rid of duplicates. After the

getting rid of the duplicates, I could use and an INSERT to populate the tables with referees and teams. Here is an example of the code I used to get rid of duplicates:

```
#gets rid of duplicates
for input in refList:
    if input not in refListU:
        t2 += 1
        refListU.append(input)
        print input[1] + " not duplicate"
```

The game table was fairly simple I just had retrieved the teamID for the away and home teams. The stats table was even simpler because I generated the gameID the same in both the game and stats table. This made the stats table almost a one by one upload from the source table. See script 4 in the Appendix for the how the tables were populated.

## Analysis:

Now that the data has been normalized in a way that it could be used in an application I did an analysis of the data. The first SQL query I wrote was to find out how many games each team had played in England's first division. The query I used to find this was:

```
1 |SELECT t.team, Count(result.team) as Games FROM (
2 |    SELECT awayteamID as team, TgameID
3 |    FROM game
4 |    UNION
5 |    SELECT hometeamID as team, TgameID
6 |    FROM game
7 |    GROUP BY TgameID
8 |
9 |) result
10 | JOIN teams t ON t.TteamID = result.team
11 | group BY t.team
12 | ORDER BY Games ASC
```

Top 10 results

Bottom 10 results

team	Games ▾	team	Games ▲
Chelsea	972	Middlesboro	38
Everton	972	Blackpool	38
Man United	972	Barnsley	38
Arsenal	972	Swindon	42
Liverpool	972	Oldham	42
Tottenham	972	Huddersfield	52
Newcastle	896	Brighton	52
Aston Villa	882	Cardiff	52
West Ham	858	Bradford	76
Man City	782	Sheffield United	80

This shows that since 1993 six teams have played in 972 games, while on the other side 5 teams have only played in one season.

Since not every team has played the same amount of games, it does not make sense to look at total goals and shots instead my query looks at the average of goals and shots to compare them. To do this I had to write two queries:

```

1 SELECT t.team, AVG(result.goals) FROM (
2   SELECT awayteamID as team, s.awayFinalGoals as
   goals, TgameID
3   FROM `game`
4   JOIN `stats` s ON s.gameID = TgameID
5   UNION
6   SELECT hometeamID as team, s.homeFinalGoals as
   goals, TgameID
7   FROM `game`
8   JOIN `stats` s ON s.gameID = TgameID
9   GROUP BY TgameID
10 ) result
11 JOIN `teams` t ON t.TteamID = result.team
12 group BY t.team
13 ORDER BY AVG(result.goals) DESC

```

```

1 SELECT t.team, AVG(result.shots) FROM (
2   SELECT awayteamID as team, s.awayShots as shots,
   TgameID
3   FROM `game`
4   JOIN `stats` s ON s.gameID = TgameID
5   UNION
6   SELECT hometeamID as team, s.homeShots as shots,
   TgameID
7   FROM `game`
8   JOIN `stats` s ON s.gameID = TgameID
9   GROUP BY TgameID
10 ) result
11 JOIN `teams` t ON t.TteamID = result.team
12 group BY t.team
13 ORDER BY AVG(result.shots) DESC

```

Top 10 Average goals

team	AVG(result.goals) ▾
Man United	1.9331
Arsenal	1.8148
Chelsea	1.7346
Liverpool	1.6975
Man City	1.6189
Tottenham	1.4866
Blackpool	1.4474
Leeds	1.3709
Newcastle	1.3594
Blackburn	1.3135

Top 10 Average shots

team	AVG(result.shots) ▾
Liverpool	14.9914
Chelsea	14.8954
Man United	14.5688
Arsenal	14.2206
Man City	13.8470
Tottenham	13.6734
QPR	13.0877
Southampton	12.2847
Bournemouth	12.0547
Everton	11.8481

Interestingly the top 6 results have the same teams in the but the next 4 results in the top 10 do not contain any of the same team. A couple things to note are the Liverpool averages the most shots per game, however Liverpool average only the fourth most goals per game. The other thing of note is that only two teams match their position of both queries: Man City and Tottenham. While looking at my results I noticed something interesting while looking at the lowest averages for goals:

team	AVG(result.goals) ▲
Huddersfield	0.7115
Cardiff	0.8654
Bradford	0.8947
Sheffield United	0.9250
Burnley	0.9518
Hull	0.9526
Brighton	0.9615
Barnsley	0.9737
Oldham	1.0000
Watford	1.0049

Oddly in the bottom ten only two of the five teams that only played one season appeared, Barnsley, and Oldham. I would of expected all 5 teams to be there.

The Next thing I wanted to look at was how dirty some teams have been. Once again because not every team has played the same amount of games so totals do not give a very good view. The first query I looked at was on average how many yellow cards did a team receive:

```

1 SELECT t.team, AVG(result.YellowCards) FROM (
2     SELECT awayteamID as team, s.awayYellowCards as
3     YellowCards, TgameID
4     FROM `game`
5     JOIN `stats` s ON s.gameID = TgameID
6     UNION
7     SELECT hometeamID as team, s.homeYellowCards as
8     YellowCards, TgameID
9     FROM `game`
10    JOIN `stats` s ON s.gameID = TgameID
11    GROUP BY TgameID
12 ) result
13 JOIN `teams` t ON t.TteamID = result.team
14 GROUP BY t.team
15 ORDER BY AVG(result.YellowCards) DESC

```

Top 3 results

team	AVG(result.YellowCards) ▾
Coventry	2.0000
Leeds	1.8816
Derby	1.8509

Looking at this result Coventry is by far the dirtiest playing team averaging 2 yellow cards a game. I gets even more compelling when looking at the average red cards per game:



team	AVG(result.YellowCards)	AVG(result.RedCards) ▾
Middlesboro	1.3158	0.1316
QPR	1.6754	0.1316
Hull	1.7368	0.1263
Blackburn	1.8014	0.1196
Leicester	1.5500	0.1107
Coventry	2.0000	0.1053

This query is ordered by the most red cards on average. On average about every 10 games Coventry receives a red card. Looking at the other top results Blackburn is the only other team that could give Coventry a run for its money. However, based on this data Coventry appears to be on average the most cautioned team in the first division.

Now soccer is a game so there is a winner and loser, sometimes there is a tie. So I looked at who has had the most wins, loses, and ties since 1993. This took three queries because of problems with UNION.

Wins:

```

1 SELECT t.team, COUNT(result.wins) FROM (
2   SELECT awayteamID as team, TgameID, s.result as wins
3   FROM game
4   JOIN stats s ON s.gameID = TgameID
5   WHERE s.result = 'A'
6   UNION
7   SELECT hometeamID as team, TgameID, s.result as wins
8   FROM game
9   JOIN stats s ON s.gameID = TgameID
10  WHERE s.result = 'H'
11  GROUP BY TgameID
12  ) result
13 JOIN teams t ON t.TteamID = result.team
14 GROUP BY t.team
15 ORDER BY COUNT(result.wins) DESC

```

team	COUNT(result.wins) ▾
Man United	611
Arsenal	538
Chelsea	532
Liverpool	494

For wins I included the top four to show how far ahead Man United is from the rest of the field.

Loses:

```

1 SELECT t.team, COUNT(result.loses) FROM (
2   SELECT awayteamID as team, TgameID, s.result as
   loses
3   FROM `game`
4   JOIN `stats` s ON s.gameID = TgameID
5   WHERE s.result = 'H'
6   UNION
7   SELECT hometeamID as team, TgameID, s.result as
   loses
8   FROM `game`
9   JOIN `stats` s ON s.gameID = TgameID
10  WHERE s.result = 'A'
11  GROUP BY TgameID
12 ) result
13 JOIN `teams` t ON t.TteamID = result.team
14 group BY t.team
15 ORDER BY COUNT(result.loses) DESC

```

team	COUNT(result.loses) ▾
West Ham	355

Ties:

```

1 SELECT t.team, COUNT(result.ties) FROM (
2   SELECT awayteamID as team, TgameID, s.result as ties
3   FROM `game`
4   JOIN `stats` s ON s.gameID = TgameID
5   WHERE s.result = 'D'
6   UNION
7   SELECT hometeamID as team, TgameID, s.result as ties
8   FROM `game`
9   JOIN `stats` s ON s.gameID = TgameID
10  WHERE s.result = 'D'
11  GROUP BY TgameID
12 ) result
13 JOIN `teams` t ON t.TteamID = result.team
14 group BY t.team
15 ORDER BY COUNT(result.ties) DESC

```

team	COUNT(result.ties) ▾
Everton	283

This is just to the top team in each category. All three teams in each category have been around the English first division a long time with West Ham the only team of the three missing in any of the seasons since 1993.

Finally, the last analysis that I did was checking what referee officiated the most home wins. Though, the referee data is far from complete with multiple seasons having the value for referee NULL. The query I used for this analysis was:

```

1 SELECT r.Original, COUNT(*)
2 FROM game
3 JOIN stats s ON s.gameID = TgameID
4 Left JOIN referee r ON s.refID = r.TrefID
5 WHERE s.result = 'H' AND s.refID IS NOT NULL
6 GROUP BY r.TrefID
7 ORDER BY COUNT(*) DESC

```

Referee with most Home wins

Original	COUNT(*)
M Dean	189

After looking at this result I realized that there is not enough context to see if there is any significance to the result. To try and get more context I ran the query for ties and losses as well.

Referee with most ties

Original	COUNT(*)
M Dean	126

Referee with most Away wins

Original	COUNT(*)
M Dean	118

After running the three queries and seeing all three have the same referee I can only concluded that M Dean has just officiated more games than everyone else.

**Conclusion:**

This project took a significant amount of work for me to complete. Though I believe I did the most I could given the time I was able to commit to it. It also taught me some very beneficial skills from building a database from scratch to being able to web scrape hundreds of CSVs automatically from a website. Finally because of the way I normalize my database I had to write some fairly complicated SQL queries to complete my analysis. All in all I enjoyed doing this project because this is similar to work I would like to do after graduating.

I have attached a file zipped with all the code I used to complete this project.

## APPENDIX:

```

1 import csv
2 import re
3 import pymysql
4 import glob
5
6 tablename = 'soccerDataSource'
7
8 conn = pymysql.connect(host='mysql.clarksonmsda.org', port=3306, user='nahomjd', passwd='X7phenom5!', db='nahomjd_is426', autocommit=True) #setup our credentials
9 cur = conn.cursor(pymysql.cursors.DictCursor)
10
11 inputs = [] # etc
12 count = 0
13
14 for files in glob.glob('C:/Users/User/Desktop/FinalProject/backupCSVs/*E0.csv'):
15     inputs.append(files)
16     print files
17
18 #determine field names
19 fieldnames = []
20 for filename in inputs:
21     with open(filename, "r") as f_in:
22         reader = csv.reader(f_in)
23         headers = next(reader)
24         for header in headers:
25             if header not in fieldnames:
26                 if len(header) > 0:
27                     fieldnames.append(header)
28                     count += 1
29                     print header
30
31 sql = "CREATE TABLE IF NOT EXISTS '"+tablename+"' ("
32 i=0
33 c2m = {}
34 m2c = {}
35
36 for field in fieldnames:
37     fn = field.lower().replace(" ", "_")
38     fn = field.lower().replace("'", "")
39
40     c2m[field] = fn
41     m2c[fn] = field
42     if i > 0:
43         sql += ",\n"
44         sql += "'"+fn+"' + " VARCHAR(255) NULL"
45         i+=1
46     sql += " "
47 sql += " );"
48
49 cur.execute(sql)
50
51 conn.close()

```

createSourceTable.py

```
3 def format(date):
4     count = 0
5     if date == "":
6         return ""
7     day = date.split("/")[0]
8
9     if len(day) == 1:
10         day = "0" + day
11
12     month = date.split("/")[1]
13
14     if len(month) == 1:
15         month = "0" + month
16
17     year = date.split("/")[2]
18
19     if len(year) == 2:
20         if int(year) > 50:
21             year = "19" + year
22         else:
23             year = "20" + year
24
25     formatted = year + "-" + month + "-" + day
26     count += 1
27
28     return formatted
```

formatDate Function

```

3 import json
4 import urllib2
5 from readCSV import upload
6
7 test = 0
8 database = "soccerDataSource"
9 #http://www.football-data.co.uk/new/ARG.csv
10 count = 0
11 base = requests.get("http://www.football-data.co.uk/data.php")
12 #glob function
13 rows = base.text.split("\n")
14 #http://www.football-data.co.uk/mmz4281/1819/E0.csv
15 baseURL = 'http://www.football-data.co.uk/'
16 for row in rows:
17     if '<tr><td>' in row:
18         country = row.split("HREF=\"")[1].split("\">")[0]
19         #print country
20
21         r = requests.get("http://www.football-data.co.uk/" + country)
22
23         lines = r.text.split("\n")
24         for line in lines:
25             #print line
26             if '<IMG SRC="Excel' in line:
27                 field = line.split("HREF=\"")[1].split("\">")[0]
28
29                 link = baseURL + field
30                 csv_ = field.split("/", 1)[1]
31                 file = csv_.replace("/", "_")
32                 print file
33
34                 data = requests.get(link)
35
36
37
38
39
40             #print file
41             f = open(file, 'wb')
42             f.write(data.text.encode('ascii', 'replace'))
43             f.close
44
45             print file + " csv created!"
46             count += 1
47
48
49
50         if '.csv"><IMG' in line:
51             field = line.split("HREF=\"")[1].split("\">")[0]
52
53             link = baseURL + field
54             file = field.split("/", 1)[1]
55
56             print file
57             data = requests.get(link)
58
59             f = open(file, 'wb')
60             f.write(data.text.encode('ascii', 'replace'))
61             f.close
62
63             print file + " csv created!"
64             count += 1
65
66 print count
67

```

downloadCSV.py



```

98 if game:
99     hTeamID = 0
100     aTeamID = 0
101
102     for row in curall:
103         sql = 'SELECT * FROM `teams` WHERE team = %s'
104         #getting hometeam ID
105         curteam = conn.cursor(pymysql.cursors.DictCursor)
106         curteam.execute(sql,(row['hometeam']))
107         for team in curteam:
108             htid = team['teamID']
109
110         sql = 'SELECT * FROM `teams` WHERE team = %s'
111         #getting awayteam ID
112         curteam = conn.cursor(pymysql.cursors.DictCursor)
113         curteam.execute(sql,(row['awayteam']))
114         for team in curteam:
115             atid = team['teamID']
116         curwareInsert = conn.cursor(pymysql.cursors.DictCursor)
117
118         sql = 'INSERT INTO game (`date`,`division`,`homeTeamID`,`awayTeamID`) VALUES \
119             (%s,%s,%s,%s);'
120         curwareInsert.execute(sql,(row['date'],row['div'],htid, atid))
121         print "intserted"
122

```

gameUpload

```

27 if ref:
28     for row in curall:
29         if row['referee']:
30             count = row['referee'].count(' ')
31             #dealing with last name first
32             if ',' in row['referee']:
33                 refList.append([row['referee'].split(" ")[1], row['referee'].split(",")[0], row['referee']])
34
35             else:
36                 #if len(row['referee'].split(" ")[0]) > 1:
37                 #print row['referee'].split()
38                 #dealing with middle intial
39                 if len(row['referee'].split()) == 3:
40                     refList.append([row['referee'].split(" ")[0], row['referee'].split(" ")[2],row['referee']])
41                 elif len(row['referee'].split()) == 2:
42                     refList.append([row['referee'].split(" ")[0], row['referee'].split(" ")[1],row['referee']])
43                 else:
44                     print "error"
45             t += 1
46
47         #gets rid of duplicates
48         for input in refList:
49             if input not in refListU:
50                 t2 += 1
51                 refListU.append(input)
52                 print input[1] + " not duplicate"
53
54         print t, t2
55
56         #print refListU
57         #insert List
58         for name in refListU:
59             if len(name) != 0:
60                 #print name
61                 print name[0] + " " + name[1] + " " + name[2]
62                 curwareInsert = conn.cursor(pymysql.cursors.DictCursor)
63
64                 sql = 'INSERT INTO referee (`fName`,`lName`, `Original`) VALUES \
65                     (%s,%s,%s);'
66                 curwareInsert.execute(sql,(name[0],name[1],name[2]))
67                 print "intserted"
68
69

```

refereesUpload

statsUpload

```

71 if team:
72     for row in curall:
73         if row['hometeam']:
74             #print row['hometeam']
75             teamList.append(row['hometeam'])
76         if row['hometeam']:
77             #print row['hometeam']
78             teamList.append(row['awayteam'])
79         #gets rid of duplicates
80         for input in teamList:
81             if input not in teamListU:
82                 t2 += 1
83                 teamListU.append(input)
84                 print input + " not duplicate"
85         #upload teams
86         for name in teamListU:
87             if len(name) != 0:
88                 #print name
89                 print name
90                 curwareInsert = conn.cursor(pymysql.cursors.DictCursor)
91
92                 sql = 'INSERT INTO teams (`team`) VALUES \
93                     (%s);'
94                 curwareInsert.execute(sql,(name))
95                 print "intserted"
96

```

teamsUpload



```

7 def upload(filename, tablename):
8     count = 0
9     conn = pymysql.connect(host='mysql.clarksonmsda.org', port=3306, user='nahomjd', passwd='X7phenom5!', db='nahomjd_is426', autocommit=True) #setup our credentials
10    cur = conn.cursor(pymysql.cursors.DictCursor)
11    with open(filename) as f:
12        data = [(k: str(v) for k, v in row.items())
13                for row in csv.DictReader(f, skipinitialspace=True)]
14
15    fields = data[0].keys()
16    date = ""
17    sql = "CREATE TABLE IF NOT EXISTS `soccerDataSource` ("
18
19    failure = 0
20    i=0
21    c2m = {}
22    m2c = {}
23    for field in fields:
24        if field == "":
25            failure += 1
26        else:
27            fn = field.lower().replace(" ", "_")
28            fn = field.lower().replace(".", "_")
29
30            c2m[field] = fn
31            m2c[fn] = field
32            if i > 0:
33                sql += ",\n"
34                sql += " " + fn + " " + " VARCHAR(255) NULL"
35                i += 1
36    sql += " ) ;"
37    cur.execute(sql)
38
39    f = open("mytable.sql", 'w')
40    f.write(sql)
41    f.close()
42
43    f1 = ','.join(m2c.keys())
44    f1 = '`' + f1 + '`'
45    t1 = ''
46    t1 = str('%s,' * len(m2c.keys()))[0:-1]
47    sql = 'INSERT INTO `'+tablename+'` ('+f1+') \n'
48    sql = 'VALUES ('+t1+')';
49
50    for row in data:
51        data_ok = False
52        fd1 = []
53        for fn in m2c.keys():
54            if fn == 'date':
55                date = format(row[m2c[fn]])
56                if len(date) > 0:
57                    data_ok = True
58                    fd1.append(date)
59                    count += 1
60            else:
61                fd1.append(row[m2c[fn]])
62            #print sql
63            #print fd1
64        if data_ok:
65            #print sql
66            cur.execute(sql, fd1)
67    print count
68    cur.close()
69    conn.close()
70    #print "failed columns" + str(failure)

```

## Upload Function

```

1  import requests
2  import csv
3  import json
4  import urllib2
5  import glob
6  from readCSV import upload
7  import pymysql
8
9  conn = pymysql.connect(host='mysql.clarksonmsda.org', port=3306, user='nahomjd', passwd='X7phenom5!', db='nahomjd_is426', autocommit=True) #setup our credentials
10  curall = conn.cursor(pymysql.cursors.DictCursor)
11
12  database = "soccerDataSource"
13  trun = conn.cursor(pymysql.cursors.DictCursor)
14  sql = 'TRUNCATE TABLE `soccerDataSource`;';
15  trun.execute(sql)
16
17  for files in glob.glob('C:/Users/User/Desktop/FinalProject/backupCSVs/*_E0.csv'):
18
19      field = files.split('\\')[1]
20      print field
21      upload(files, database)
22      print "uploaded"

```

uploadSource.py