

```

01. import os
02.
03. from calendar import monthrange
04. from nba_py import game, Scoreboard
05. import pandas as p
06. import argparse
07.
08.
09. MONTHS_BEFORE_NEWYEAR = range(10,13)
10. MONTHS_AFTER_NEWYEAR = range(1,5)
11. UNNEEDED_FEATURES = ['GAME_ID', 'TEAM_ID', 'TEAM_ABBREVIATION', 'TEAM_CITY', 'MIN', 'FGA', 'FGM',
12. NO_DEFICIT_FEATURES = ['PCT', 'TEAM', 'ID', 'MIN', 'OUTCOME', 'PTS', 'PLUS_MINUS']
13. AWAY_WIN = [0,1]
14. HOME_WIN = [1,0]
15. CSV = '.csv'
16.
17.
18.
19. class Create_csv_files(object):
20.     """
21.     Creating csv file containing the data of an NBA season
22.     """
23.
24.
25.     def __init__(self, root_dir, season):
26.         """
27.         Initialize the object
28.         @param root_dir: a path to the root directory which all the files will be stored
29.         @type root_dir: string
30.         @param season: the years the season is being played
31.         @type season: list of strings, containing 2 strings
32.         """
33.         self._root_dir = root_dir + '/'
34.         # Season must be a list with two strings in the form of [2012,2013]
35.         if not (len(season) == 2 and (season[1] - season[0] == 1)):
36.             print('season in not correct')
37.         self._season = season
38.         self.df = p.DataFrame()
39.         self._current_teams = []
40.
41.     def _create_directory(self, curr_dir):
42.         """
43.         Create a directory to contain the new csv files
44.         """
45.         if not os.path.exists(curr_dir):
46.             os.makedirs(curr_dir)
47.
48.     def season_name(self):
49.         """
50.         @return: the season name as a form of string
51.         @rtype: str
52.         """
53.         return str(self._season[0]) + '-' + str(self._season[1])
54.
55.     def scrape_season(self, test=False):
56.         """
57.         Scrape a full season from the nba stats website
58.         @param test: for debugging
59.         @type test: bool
60.         """
61.         season_dir = self._root_dir + '/' + self.season_name() + '/'
62.         month_years = [(self._season[0], month) for month in MONTHS_BEFORE_NEWYEAR] + \
63.             [(self._season[1], month) for month in MONTHS_AFTER_NEWYEAR]
64.
65.         self._create_directory(season_dir)
66.
67.         if test:
68.             # For debugging
69.             self.scrape_month(season_dir, 2014, 11, first_day=22,last_day=25)
70.         else:
71.             # Scrape the all season
72.             for year, month in month_years:
73.                 print('currently scraping: {}/{}'.format(year,month))

```

```

74.         self.scrape_month(season_dir, year, month)
75.
76.         all_boxscores_file = season_dir + 'all.csv'
77.         print("create a csv file with all the games named ", all_boxscores_file)
78.         self.df.to_csv(all_boxscores_file)
79.
80.     def scrape_month(self, root_dir, year, month, first_day=1, last_day=31):
81.         """
82.         Create a directory containing all the boxscore of the entered month
83.         @param root_dir: the directory to hold the current month data
84.         @type root_dir: str
85.         @param year: the current year to gather data from
86.         @type year: int
87.         @param month: the current month to gather data from
88.         @type month: int
89.         @param first_day: the first day of games in the month
90.         @type first_day: int
91.         @param last_day: the last day of games in the month
92.         @type last_day: int
93.         """
94.         last_day = min([monthrange(year, month)[1], last_day])
95.         curr_dir = root_dir + str(month) + '/'
96.
97.         self._create_directory(curr_dir)
98.
99.         for day in range(first_day, last_day+1):
100.             print("currently scraping {day}/{month}/{year}".format(day=day, month=month, year=year))
101.             self.process_scoreboard(year, month, day, curr_dir)
102.
103.     def _fill_teams(self, sb):
104.         """
105.         Gather all the teams of the current season
106.         @param sb: data about the current day
107.         @type sb: ScoreBoard Object
108.         """
109.         es = sb.east_conf_standings_by_day()
110.         ws = sb.west_conf_standings_by_day()
111.
112.         self._current_teams += list(ws['TEAM']) + list(es['TEAM'])
113.
114.     def process_scoreboard(self, root_dir, year, month, day ):
115.         """
116.         Create csv files with game id inside a directory by day
117.         @param root_dir: the directory to hold the current month data
118.         @type root_dir: str
119.         @param year: the current year to gather data from
120.         @type year: int
121.         @param month: the current month to gather data from
122.         @type month: int
123.         @param day: the current day of games
124.         @type day: int
125.         """
126.         sb = Scoreboard(month=month, day=day, year=year)
127.         curr_dir = root_dir + str(day) + '/'
128.
129.         if not self._current_teams:
130.             self._fill_teams(sb)
131.         self._create_directory(curr_dir)
132.
133.         for gid in sb.available().GAME_ID:
134.             self.boxscore_to_csv(gid, curr_dir)
135.
136.     def remove_columns(self, ts):
137.         """
138.         Remove unnecessary features we don't need to be in the final data
139.         """
140.         for feature in UNNEEDED_FEATURES:
141.             del ts[feature]
142.         return ts
143.
144.     def add_deficit_column(self, ts, feature):
145.         """
146.         Add data containing the deficit between the teams
147.         """

```

```

148.         ts['DEFICIT_'+feature] = ts[feature] / sum(ts[feature])
149.     return ts
150.
151.     def add_columns(self, ts):
152.         """
153.         Adding columns to a given dataframe
154.         @param ts: box score
155.         @type ts: dataframe
156.         """
157.         # Adding 2 pt field goals data
158.         ts['FG2M'] = ts['FGM'] - ts['FG3M']
159.         ts['FG2A'] = ts['FGA'] - ts['FG3A']
160.         ts['FG2_PCT'] = ts['FG2M'] / ts['FG2A']
161.
162.         # 1 represents win 0 represents loss
163.
164.         ts['OUTCOME'] = HOME_WIN if ts['PLUS_MINUS'][0] > 0 else AWAY_WIN
165.
166.         # Add deficit
167.         for c in ts.columns:
168.             if any(s in c for s in NO_DEFICIT_FEATURES):
169.                 continue
170.             ts = self.add_deficit_column(ts, c)
171.
172.         return ts
173.
174.     def irrelevant_data(self, ts):
175.         """
176.         Test if the current receive data needs to be deleted
177.         """
178.         if len(ts.index) == 0:
179.             return True
180.
181.         if len(ts.dropna()) < 2:
182.             return True
183.
184.         if not all(team in self._current_teams for team in ts['TEAM_CITY']):
185.             return True
186.
187.         return False
188.
189.
190.     # TODO: add manipulation over the data
191.     def manipulate_df(self, ts):
192.         """
193.         Manipulate data to fit for the research
194.         """
195.         if self.irrelevant_data(ts):
196.             return None
197.
198.         ts = self.add_columns(ts)
199.
200.         ts = self.remove_columns(ts)
201.
202.         return ts
203.
204.     def boxscore_to_csv(self, gid, root_dir):
205.         """
206.         Create a csv file from a boxscore according to the game id
207.         """
208.         bs = game.Boxscore(gid)
209.         team_stats = bs.team_stats()
210.         team_stats = self.manipulate_df(team_stats)
211.         if team_stats is None:
212.             return
213.
214.         csv_name = root_dir + gid + CSV
215.         team_stats.to_csv(csv_name)
216.
217.         self.df = self.df.append(team_stats, ignore_index=True)
218.
219.
220.     def main():
221.         parser = argparse.ArgumentParser(description="")

```

```
222.     parser.add_argument('--directory', dest='dir')
223.     parser.add_argument('--season', dest='se', nargs='+', type=int)
224.     args = parser.parse_args()
225.
226.     my_tool = Create_csv_files(args.dir, args.se)
227.     my_tool.scrape_season()
228.
229.
230. if __name__ == "__main__":
231.     main()
```