

Final Project Data Check point

1. Project code:

https://github.com/cwtseng/SI507_final_proj.git

2. Data sources:

A. Origin URL: <https://www.nba.com/stats/teams/boxscores/>

i. Website data example

1662 Rows Page 31 of 34																							
TEAM	MATCH UP	GAME DATE	W/L	MIN	PTS	FGM	FGA	FG%	3PM	3PA	3P%	FTM	FTA	FT%	OREB	DREB	REB	AST	STL	BLK	TOV	PF	+/-
IND	IND vs. NYK	01/02/2021	L	240	102	34	84	40.5	19	50	38.0	15	18	83.3	5	27	32	22	6	4	9	21	-4
NYK	NYK @ IND	01/02/2021	W	240	106	42	83	50.6	12	27	44.4	10	13	76.9	10	41	51	23	4	3	16	16	4
HOU	HOU vs. SAC	01/02/2021	W	240	102	36	81	44.4	13	41	31.7	17	23	73.9	6	42	48	15	10	10	14	22	8

B. Explanation of data label: <https://ir.nba.com/how-to-read-a-box-score/>

i. Data label explanation

- MIN = Minutes
- FGM = Field-goals made
- FGA = Field-goals attempted
- FG% = Field goal percentage
- 3PM = 3-pointers made
- 3PA = 3-pointers attempted
- 3P% = 3-point percentage
- FTM = Free throws made
- FTA = Free throws attempted
- FT% = Free throw percentage
- OREB = Offensive rebounds
- DREB = Defensive rebounds
- REB = Total rebounds
- AST = Assists
- TOV = Turnovers
- STL = Steals
- BLK = Blocked shots
- PF = Personal fouls
- PTS = Points scored
- +/- = Plus/Minus

C. Raw data format: HTML

D. How to access the data?

- Grab multiple pages raw HTMLs with Selenium(webdriver and Select)
- Parse the raw HTMLs via BeautifulSoup
- Organize data with pandas dataframe
- Store parsed data in cache and database

E. Summary of data:

- >400 records available.
- >400 records retrieved.
- Retrieved all data available on the website, because all of them are useful.
- The row of data indicates different teams and different dates result
- The column of data indicates different key performance result in the team.
- The data can be processing with specified key performance(ex: field goal percent, assist, rebound), teams(ex: Lakers, Clippers), and interesting period(ex: start date and end date)

- vii. For a single team, I can plot different key performance over interesting period. Also, user can apply statistics options, ex: average, maximum, minimum, standard deviation within that period.
- viii. For multiple teams, the program can show their comparison of key performance in the same plot It can also applied the statistics result as well

F. Evidence of caching:

- i. Example of cache file 2020-21_RegularSeason.json, store all the crawling data

```
{
  "columns": [
    "TEAM_NAME", "MATCHUP", "gdate", "WL", "MIN", "PTS", "FGM", "FGA", "FG_PCT", "FG3M", "FG3A", "FG3_PCT", "FTM", "FTA", "FT_PCT", "OREB", "DREB", "REB", "AST", "STL", "BLK", "TOV", "PF", "PLUS_MINUS", "data": [
      ["GSW", "GSW vs. HOU", "2021-04-10", "W", "240,125,49,91,0.538,14,35,0.4,13,16,0.813,9,36,45,31,13,5,16,19,16], ["HOU", "HOU @ GSW", "2021-04-10", "L", "240,109,41,88,0.466,13,37,0.351,14,20,0.7,13,32,45,25,8,3,18,17, -16], ["POR", "POR vs. DET", "2021-04-10", "W", "240,118,43,90,0.478,13,29,0.448,19,24,0.792,14,35,49,18,5,8,16,20,15], ["DET", "DET @ POR", "2021-04-10", "L", "240,103,37,83,0.446,9,26,0.346,20,28,0.714,9,29,38,25,9,6,12,23, -15], ["WAS", "WAS @ PHX", "2021-04-10", "L", "240,106,39,84,0.464,9,25,0.36,19,28,0.679,7,42,49,27,2,1,16,12, -28], ["PHX", "PHX vs. WAS", "2021-04-10", "W", "240,134,55,107,0.514,17,43,0.395,7,7,1,0,5,40,45,36,8,3,3,18,28], ["SAC", "SAC @ UTA", "2021-04-10", "L", "240,112,42,88,0.477,14,42,0.333,14,22,0.636,4,34,38,23,8,4,10,23, -16], ["UTA", "UTA vs. SAC", "2021-04-10", "W", "240,128,40,88,0.455,18,50,0.36,30,35,0.857,11,41,52,22,8,5,10,17,16], ["OKC", "OKC vs. PHI", "2021-04-10", "L", "240,93,38,76,0.5,8,20,0.4,9,19,0.474,6,39,45,21,3,4,23,21, -24], ["PHI", "PHI @ OKC", "2021-04-10", "W", "240,117,44,94,0.468,12,38,0.316,17,25,0.68,13,34,47,24,11,10,10,17,24], ["LAL", "LAL @ BKN", "2021-04-10", "W", "240,126,47,93,0.505,19,34,0.559,13,20,0.65,12,35,47,30,12,4,15,24,25], ["BKN", "BKN vs. LAL", "2021-04-10", "L", "240,101,35,80,0.438,5,27,0.185,26,30,0.867,9,33,42,19,6,10,19,20, -25], ["CLE", "CLE vs. TOR", "2021-04-10", "L", "240,115,39,86,0.453,15,35,0.429,22,26,0.846,7,25,32,27,11,5,17,13, -20], ["TOR", "TOR @ CLE", "2021-04-10", "W", "240,135,53,86,0.616,17,32,0.531,12,14,0.857,5,34,39,29,12,11,19,18,20], ["GSW", "GSW vs. WAS", "2021-04-09", "L", "240,107,41,88,0.466,11,36,0.306,14,21,0.667,8,33,41,27,15,3,13,20, -3], ["WAS", "WAS @ GSW", "2021-04-09", "W", "240,110,43,89,0.483,5,19,0.263,19,24,0.792,11,36,47,30,6,1,18,16,3], ["LAC", "LAC vs. HOU", "2021-04-09", "W", "240,126,48,88,0.545,19,37,0.514,11,13,0.846,4,34,38,31,10,4,10,24,17], ["HOU", "HOU @ LAC", "2021-04-09", "L", "240,109,39,83,0.47,12,37,0.324,19,27,0.704,12,29,41,28,3,6,14,15, -17], ["SAS", "SAS @ DEN", "2021-04-09", "L", "240,119,43,93,0.462,9,27,0.333,24,28,0.857,11,27,38,33,10,2,9,21, -2], ["DEN", "DEN vs. SAS", "2021-04-09", "W", "240,121,43,79,0.544,14,27,0.519,21,23,0.913,9,37,46,32,5,4,18,22,2], ["CHA", "CHA @ MIL", "2021-04-09", "W", "240,127,44,91,0.484,19,50,0.38,20,22,0.909,13,36,49,35,4,10,20,19,8], ["MIL", "MIL vs. CHA", "2021-04-09", "L", "240,119,44,98,0.449,14,38,0.368,17,23,0.739,13,35,48,23,10,3,17,14, -8], ["PHI", "PHI @ NOP", "2021-04-09", "L", "240,94,33,79,0.418,11,31,0.355,17,24,0.708,4,43,47,20,8,5,19,23, -7], ["NOP", "NOP vs. PHI", "2021-04-09", "W", "240,101,39,90,0.433,4,22,0.182,19,30,0.633,10,46,56,22,11,7,14,17,7], ["NYK", "NYK vs. MEM", "2021-04-09", "W", "265,133,45,96,0.469,14,32,0.438,29,33,0.879,11,37,48,24,7,7,14,25,4], ["MEM", "MEM @ NYK", "2021-04-09", "L", "265,129,47,94,0.5,14,30,0.467,21,34,0.618,10,33,43,33,7,6,14,27, -4], ["BOS", "BOS vs. MIN", "2021-04-09", "W", "265,145,47,87,0.54,19,42,0.452,32,34,0.941,9,36,45,32,9,5,21,25,9], ["MIN", "MIN @ BOS", "2021-04-09", "L", "265,145,47,87,0.54,19,42,0.452,32,34,0.941,9,36,45,32,9,5,21,25,9]
    ]
  ]
}
```

3. Database:

- A. Database has been used, single table with all team and performance result. The different Id is showing the different date.
- B. Only primary key is used in the database
- C. Part of snapshot of database

Table: teams		Filter in any column																					
	Id	TEAM_NAME	MATCHUP	gdate	WL	MIN	PTS	FGM	FGA	FG_PCT	FG3M	FG3A	FG3_PCT	FTM	FTA	FT_PCT	OREB	DREB	REB	AST	STL	BLK	TOV
	過濾	過濾	過濾	過濾	過濾	過濾	過濾	過濾	過濾	過濾	過濾	過濾	過濾	過濾	過濾	過濾	過濾	過濾	過濾	過濾	過濾	過濾	過濾
1	1	MIL	MIL @ ATL	2021-04-15	W	240	120	46	95	0.48	16	43	0.37	12	17	0.71	14	37	51	25	8	1	9
2	2	ATL	ATL vs. MIL	2021-04-15	L	240	109	39	86	0.45	14	34	0.41	17	19	0.90	8	31	39	23	5	4	12
3	3	LAL	LAL vs. BOS	2021-04-15	L	240	113	44	98	0.45	12	37	0.32	13	19	0.68	11	25	36	28	12	2	8
4	4	BOS	BOS @ LAL	2021-04-15	W	240	121	48	85	0.57	14	32	0.44	11	13	0.85	7	41	48	33	4	6	21
5	5	SAC	SAC @ PHX	2021-04-15	L	240	114	43	78	0.55	14	30	0.47	14	19	0.74	4	28	32	22	7	4	15
6	6	PHX	PHX vs. SAC	2021-04-15	W	240	122	48	86	0.56	14	38	0.37	12	14	0.86	10	29	39	30	9	4	11
7	7	GSW	GSW @ CLE	2021-04-15	W	240	119	47	92	0.51	17	46	0.37	8	10	0.80	9	38	47	33	9	5	16
8	8	CLE	CLE vs. GSW	2021-04-15	L	240	101	35	81	0.43	10	31	0.32	21	27	0.78	8	31	39	21	11	2	15
9	9	TOR	TOR vs. SAS	2021-04-14	W	240	117	41	88	0.47	14	30	0.47	21	25	0.84	13	41	54	21	8	5	16
10	10	SAS	SAS @ TOR	2021-04-14	L	240	112	38	89	0.43	17	39	0.44	19	27	0.70	12	29	41	26	6	7	13
11	11	MEM	MEM vs. DAL	2021-04-14	L	240	113	42	90	0.47	13	34	0.38	16	21	0.76	11	34	45	27	7	3	7
12	12	DAL	DAL @ MEM	2021-04-14	W	240	114	44	88	0.50	13	40	0.33	13	17	0.77	8	37	45	25	4	3	11

4. Interaction and Presentation:

- A. First input message box would ask you to input team names, abbreviations same as official NBA website. Ex: TOR LAC DET ...
- B. Second input message box would ask you to input options:
 - i. [MIN| PTS|FGM| ...]: the available numerical key performance you are going to display
 - ii. [none|avg|max|min|std]: statistic options for the data within interesting period
 - iii. <str1> : starting date of capture, format should be ex:'2021-04-15'

- iv. <str2>: end date of capture, format should be ex: '2021-04-15'
- v. <integer>: how many page number for program to crawl
- vi. [none|plot]: whether to plot data or just show table
- vii. [none|cache|db]: whether to use cache or db or direct fetch from website