



# Extract Transform Load NBA data

Project 2 - ETL

Jasmine Jones, Nicholas Noark, Haydn Whitmyer

# Extract: API-NBA

- API-NBA - “freemium” API from rapidapi.com
- Limited to 10 calls/minute and 100 calls/day
- Contains data on NBA Stats DATA, Games, Livescore, Standings, Statistics, Teams, Players, Seasons, Leagues.
- Updated on a semi-regular basis, has info on recent games (up to 10/16/21)
- Selected to pull from Game Details, Team (ID), and Players (ID)

	game_id	season_year	arena	city	country	start_time_UTC	game_duration	home_team	home_score	home_leader_id	away_team	away_score	away_leader_id
0	9061	2020	United Center	Chicago	USA	2021-04-15T00:00:00.000Z	2:07	Chicago Bulls	106	534	Orlando Magic	115	160
1	1067	2015				2016-03-10T02:30:00.000Z	2:18	Oklahoma City Thunder	120	153	LA Clippers	108	207
2	5906	2018	Barclays Center	Brooklyn	USA	2019-03-30T22:00:00.000Z	2:08	Brooklyn Nets	110	462	Boston Celtics	96	227
3	4431	2018	Vivint Smart Home Arena	Salt Lake City	USA	2018-10-23T01:00:00.000Z	2:25	Utah Jazz	84	121	Memphis Grizzlies	92	114
4	3000	2017	American Airlines Center	Dallas	USA	2017-10-29T00:30:00.000Z	2:08	Dallas Mavericks	110	36	Philadelphia 76ers	112	159
5	420	2015	Target Center	Minneapolis		2015-12-08T01:00:00.000Z	2:39	Minnesota Timberwolves	106	308	LA Clippers	110	286
6	1642	2016	Talking Stick Resort Arena	Phoenix		2016-11-10T02:00:00.000Z	2:16	Phoenix Suns	107	59	Detroit Pistons	100	89
7	130	2015	Quicken Loans Arena	Cleveland		2015-10-30T23:00:00.000Z	2:17	Cleveland Cavaliers	102	265	Miami Heat	92	536
8	5240	2018	Bankers Life Fieldhouse	Indianapolis	USA	2019-02-14T00:00:00.000Z	2:10	Indiana Pacers	97	60	Milwaukee Bucks	106	20
9	5123	2018	Chesapeake Energy Arena	Oklahoma City	USA	2019-01-27T23:00:00.000Z	2:20	Oklahoma City Thunder	118	189	Milwaukee Bucks	112	20

```
# game details info API to call

details_base_url = "https://api-nba-v1.p.rapidapi.com/gameDetails/"

game_numbers = list(np.random.randint(10861, size=10))
game_numbers

game_info = []
game_not_found = []

for number in game_numbers:

    game_url = details_base_url + str(number)

    try:
        game_response = requests.get(game_url, headers=headers).json()

        gameId = game_response['api']['game'][0]['gameId']
        seasonYear = game_response['api']['game'][0]['seasonYear']
        arena = game_response['api']['game'][0]['arena']
        city = game_response['api']['game'][0]['city']
        country = game_response['api']['game'][0]['country']
        startTimeUTC = game_response['api']['game'][0]['startTimeUTC']
        gameDuration = game_response['api']['game'][0]['gameDuration']
        vTeam = game_response['api']['game'][0]['vTeam']['fullName']
        vTeamLeader = game_response['api']['game'][0]['vTeam']['leaders'][0]['playerId']
        awayScore = game_response['api']['game'][0]['vTeam']['score']['points']
        hTeam = game_response['api']['game'][0]['hTeam']['fullName']
        hTeamLeader = game_response['api']['game'][0]['hTeam']['leaders'][0]['playerId']
        homeScore = game_response['api']['game'][0]['hTeam']['score']['points']

        game_info.append({
            'game_id': gameId,
            'season_year': seasonYear,
            'arena': arena,
            'city': city,
            'country': country,
            'start_time_UTC': startTimeUTC,
            'game_duration': gameDuration,
            'home_team': hTeam,
            'home_score': homeScore,
            'home_leader_id': hTeamLeader,
            'away_team': vTeam,
            'away_score': awayScore,
            'away_leader_id': vTeamLeader
        })

    except:
        game_not_found.append(gameId)
        pass

# store any values not found and pass to keep the loop running
```

# Extract: Basketball-Reference & SQLite Database

- Using Pandas and BeautifulSoup, scraped per 100 possessions and per 36 minute player statistics for the 2020-21 NBA season
- From a NBA SQLite database found on kaggle.com, pulled player draft data and player salaries for the 2020-21 NBA season
- The SQLite database was found at <https://www.kaggle.com/wyattowalsh/basketball>

Player Per 100 Poss															Share & Export	When table is sorted, hide non-qualifiers for rate stats	Glossary	Hide Partial Rows									
RA	Player	Pos	Age	Tm	G	GS	MP	FG	FGA	FG%	3P	3PA	3P%	FT	FTA	FT%	ORB	TREB	AST	STL	BLK	TOV	PF	PTS	DRtg	ORtg	
1	Travis Alexander	PF	22	PHI	6	4	737	8.4	15.4	54.6	3.0	5.1	58.8	6.4	13.5	54.8	14	7.4	3.99	4.5	1.5	2.6	1.9	29	32	107	
2	Joel Embiid	PF	29	PHI	6	4	737	8.4	15.4	54.6	3.0	5.1	58.8	6.4	13.5	54.8	14	7.4	3.99	4.5	1.5	2.6	1.9	29	32	107	
3	Joel Embiid	PF	24	PHI	7	0	18	2.4	30.9	7.8	0.0	5.2	0.0	2.4	18.7	12.7	0.0	0.0	7.8	7.8	5.2	0.0	0.0	2.4	5.2	49	
4	Steven Adams	C	27	SAS	59	1605	5.6	9.2	61.0	0.1	0.0	0.4	0.1	4.0	1.8	44.4	6.4	9.0	13.4	3.3	1.4	1.2	3.3	3.4	13.1	122	
5	Dwain Howard	C	30	SAS	44	1040	10.4	18.5	57.0	0.0	0.1	0.0	0.1	2.0	0.5	40.0	4.7	8.4	10.3	4.0	1.7	1.5	3.9	1.4	27.7	122	
6	Laurie R. King	C	18	WOT	24	23	474	10.1	21.3	47.9	2.2	5.4	38.8	7.8	18.5	42.0	2.4	3.4	47.2	1.4	7.1	8.5	3.5	0.8	1.1	1.9	1.4
7	Laurie R. King	C	18	SAS	12	18	184	10.3	21.3	48.4	2.4	4.7	50.8	7.8	18.5	42.0	2.4	3.4	47.2	1.4	7.1	8.5	3.5	0.8	1.1	1.9	1.4
8	Laurie R. King	C	18	PHI	5	8	130	8.3	17.4	48.1	1.5	1.9	78.0	7.8	18.5	42.0	2.4	3.4	47.2	1.4	7.1	8.5	3.5	0.8	1.1	1.9	1.4
9	Travis Alexander	SG	22	PHI	15	0	47	3.2	15.4	20.8	2.1	5.9	35.7	2.2	1.1	33.0	1.1	2.1	51.0	2.1	8.4	16.5	4.3	0.8	1.2	2.1	6.5
10	Travis Alexander	SG	22	PHI	45	12	1007	8.1	15.8	51.9	1.5	1.5	100.0	3.4	5.2	65.4	1.5	2.1	71.7	1.5	6.2	6.8	2.2	1.0	1.2	1.2	1.2
11	Grant Hill	SG	25	PHI	52	38	1229	6.4	15.7	41.8	4.1	10.4	39.1	2.5	5.2	47.9	3.4	3.5	44.4	6.7	5.4	4.1	1.7	0.3	1.8	2.7	20.2
12	Jacques Adams	C	22	PHI	43	45	1884	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2
13	Jacques Adams	C	20	PHI	15	0	120	5.4	9.8	55.7	0.0	0.0	0.0	6.4	9.4	67.9	6.2	7.4	83.5	10.8	3.0	1.1	1.9	1.2	1.2	1.2	
14	Jacques Adams	C	19	PHI	15	0	1044	8.1	13.3	60.9	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2
15	Al Hassan Adams	PF	20	PHI	23	14	524	8.1	13.3	60.9	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2
16	Al Hassan Adams	PF	20	PHI	17	14	507	8.1	13.3	60.9	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2
17	Al Hassan Adams	PF	20	PHI	6	0	67	1.4	2.2	63.6	0.7	1.2	58.3	1.7	2.9	58.3	1.7	2.9	58.3	1.7	2.9	58.3	1.7	2.9	58.3	1.7	2.9
18	Grant Hill	SG	25	PHI	49	1887	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2	
19	Grant Hill	SG	25	PHI	43	45	1884	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2
20	Grant Hill	SG	25	PHI	15	0	56	1.6	6.7	24.0	0.0	0.0	0.0	2.4	6.7	35.8	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
21	Grant Hill	SG	25	PHI	17	3	151	1.7	11.7	14.6	0.0	0.0	0.0	2.4	11.7	20.5	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
22	Grant Hill	SG	25	PHI	49	1887	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2	
23	Grant Hill	SG	25	PHI	43	45	1884	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2
24	Grant Hill	SG	25	PHI	15	0	56	1.6	6.7	24.0	0.0	0.0	0.0	2.4	6.7	35.8	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
25	Grant Hill	SG	25	PHI	17	3	151	1.7	11.7	14.6	0.0	0.0	0.0	2.4	11.7	20.5	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
26	Grant Hill	SG	25	PHI	49	1887	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2	
27	Grant Hill	SG	25	PHI	43	45	1884	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2
28	Grant Hill	SG	25	PHI	15	0	56	1.6	6.7	24.0	0.0	0.0	0.0	2.4	6.7	35.8	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
29	Grant Hill	SG	25	PHI	17	3	151	1.7	11.7	14.6	0.0	0.0	0.0	2.4	11.7	20.5	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
30	Grant Hill	SG	25	PHI	49	1887	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2	
31	Grant Hill	SG	25	PHI	43	45	1884	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2
32	Grant Hill	SG	25	PHI	15	0	56	1.6	6.7	24.0	0.0	0.0	0.0	2.4	6.7	35.8	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
33	Grant Hill	SG	25	PHI	17	3	151	1.7	11.7	14.6	0.0	0.0	0.0	2.4	11.7	20.5	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
34	Grant Hill	SG	25	PHI	49	1887	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2	
35	Grant Hill	SG	25	PHI	43	45	1884	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2
36	Grant Hill	SG	25	PHI	15	0	56	1.6	6.7	24.0	0.0	0.0	0.0	2.4	6.7	35.8	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
37	Grant Hill	SG	25	PHI	17	3	151	1.7	11.7	14.6	0.0	0.0	0.0	2.4	11.7	20.5	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
38	Grant Hill	SG	25	PHI	49	1887	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2	
39	Grant Hill	SG	25	PHI	43	45	1884	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2
40	Grant Hill	SG	25	PHI	15	0	56	1.6	6.7	24.0	0.0	0.0	0.0	2.4	6.7	35.8	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
41	Grant Hill	SG	25	PHI	17	3	151	1.7	11.7	14.6	0.0	0.0	0.0	2.4	11.7	20.5	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
42	Grant Hill	SG	25	PHI	49	1887	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2	
43	Grant Hill	SG	25	PHI	43	45	1884	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2
44	Grant Hill	SG	25	PHI	15	0	56	1.6	6.7	24.0	0.0	0.0	0.0	2.4	6.7	35.8	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
45	Grant Hill	SG	25	PHI	17	3	151	1.7	11.7	14.6	0.0	0.0	0.0	2.4	11.7	20.5	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
46	Grant Hill	SG	25	PHI	49	1887	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2	
47	Grant Hill	SG	25	PHI	43	45	1884	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2
48	Grant Hill	SG	25	PHI	15	0	56	1.6	6.7	24.0	0.0	0.0	0.0	2.4	6.7	35.8	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
49	Grant Hill	SG	25	PHI	17	3	151	1.7	11.7	14.6	0.0	0.0	0.0	2.4	11.7	20.5	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
50	Grant Hill	SG	25	PHI	49	1887	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2	
51	Grant Hill	SG	25	PHI	43	45	1884	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2	50.6	5.4	7.6	72.0	5.2	11.5	16.4	2.8	0.8	2.4	2.4	2.2
52	Grant Hill	SG	25	PHI	15	0	56	1.6	6.7	24.0	0.0	0.0	0.0	2.4	6.7	35.8	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
53	Grant Hill	SG	25	PHI	17	3	151	1.7	11.7	14.6	0.0	0.0	0.0	2.4	11.7	20.5	2.4	2.4	60.0	1.6	6.7	35.8	2.4	2.4	60.0	1.6	6.7
54	Grant Hill	SG	25	PHI	49	1887	7.8	15.7	50.0	0.2	0.4	33.4	7.7	15.2													

# Transform

- Once extracted, data was cleaned using pandas. Statistics such as FG%, 3P%, and FT% were removed from the per 36 minutes data since rate stats stay the same across per 100 possession stats and per 36 minutes
- Inactive players filtered out of the SQLite data then merged with the draft data from kaggle.com
- Pulled relevant basic info (name, team, DOB, city, etc.) on players and teams from API; Basketball-Reference contains more complete statistics
- Many players from outside USA lacked data, these rows were removed with .dropna()
- Reset indices to teamID, playerId, and gameId using .set\_index()

```
stats_per36 = pd.DataFrame(player_stats.columns-headers)
stats_per36.head()
```

	Player	Pos	Age	Tm	GS	MP	FG	FGA	FG3	...	FT	
0	Brooklynn Acheson	PF	21	MIA	61	4	777	6.3	11.1	...	546	...
1	Jaylen Adams	PG	24	MIL	7	0	18	2.8	16.8	...	125	...
2	Steven Adams	C	27	MOP	58	1065	4.2	6.9	854	...	484	...
3	Ben Adreony	C	23	MIA	64	1643	7.7	11.4	359	...	199	...
4	LaMarcus Aldridge	C	35	TOR	26	33	674	7.5	15.3	...	457	...

  

	ORB	DRB	TRE	AST	STL	BLK	TOV	PF	PTS
0	3.0	6.8	19.2	1.4	0.8	1.4	2.1	4.4	16.8
1	0.8	6.8	8.8	4.8	0.8	0.8	0.8	2.8	4.8
2	4.8	6.8	11.5	2.5	1.2	0.9	1.7	2.5	9.8
3	2.4	7.2	9.8	5.8	1.3	1.1	2.8	2.4	26.1
4	1.8	1.3	6.3	2.8	0.8	1.5	1.4	1.5	18.4

[5 rows x 28 columns]

Cleaning up redundancies

```
stats_per36 = stats_per36[[['Player', 'Pos', 'Age', 'Tm', 'GS', 'MP', 'FG', 'FGA', 'FG3', 'FG3A', 'FT', 'FTA', 'STL', 'BLK', 'TOV', 'PF', 'PTS', 'ORB', 'DRB', 'TRE', 'AST', 'BLK', 'TOV', 'PF', 'PTS']]]
stats_per36.head()
```

```
stats_per36.to_csv('stats_per36.csv', index=False, header=True)
```

```
salaries_21 = pd.read_sql("""SELECT *
FROM Player_Salary
WHERE slugSeason = "2020-21", """, conn)
salaries_21 = salaries_21[['nameTeam', 'namePlayer', 'value']]
salaries_21 = salaries_21.rename(columns={'nameTeam': 'Team', 'namePlayer': 'full_name'})
salaries_21.head()
```

	Team	full_name	value
0	Atlanta Hawks	Bogdan Bogdanovic	18000000.0
1	Atlanta Hawks	Brandon Goodwin	1701593.0
2	Atlanta Hawks	Bruno Fernando	1517981.0
3	Atlanta Hawks	Cam Reddish	4458000.0
4	Atlanta Hawks	Clint Capela	16000000.0

Pull draft positions of each player and merge with active players

```
draft = pd.read_sql("""SELECT *
FROM Draft""", conn)
draft = draft[['yearDraft', 'numberPickOverall', 'numberRound', 'namePlayer', 'idPlayer', 'nameTeam']]
draft = draft.rename(columns={'yearDraft': 'Year_Drafted', 'numberPickOverall': 'Overall_Pick', 'numberRound': 'Round', 'namePlayer': 'full_name', 'idPlayer': 'id', 'nameTeam': 'Drafted_by'})
active_info = pd.merge(active, draft, on='full_name')
active_info.head()
```

	id_x	full_name	Year_Drafted	Overall_Pick	Round	id_y	Drafted_by
0	203500	Steven Adams	2013.0	12.0	1.0	203500.0	Oklahoma City Thunder
1	1628389	Sam Adreony	2017.0	14.0	1.0	1628389.0	Miami Heat
2	200746	LaMarcus Aldridge	2006.0	2.0	1.0	200746.0	Chicago Bulls
3	1629638	Nickel Alexander-Walker	2019.0	17.0	1.0	1629638.0	Brooklyn Nets

# PostgreSQL

- PostgreSQL database, “mybasketball”
- Used the sqlite data from Kaggle’s Basketball Dataset to create csv’s (<https://www.kaggle.com/wyattowalsh/basketball>)
- created tables related to data from csv files, once they were cleaned and scrapped PostgreSQL
- uploaded the csv files to the corresponding tables using SQLAlchemy

```
[3] import pandas as pd
import psycopg2

conn = psycopg2.connect(
    database='postgres', user='postgres', password='Kabrija01', host='127.0.0.1', port= '5432'
)
conn.autocommit = True

Python

[5] cursor = conn.cursor()
sql = '''CREATE database mybasketball'''
cursor.execute(sql)

Python

[6]
```

```
CREATE TABLE active_info (
    id INT,
    full_name VARCHAR UNIQUE NOT NULL PRIMARY KEY,
    year_drafted DECIMAL,
    overall DECIMAL,
    pick DECIMAL,
    round_ DECIMAL,
    id_y DECIMAL,
    drafted_by VARCHAR
);
```

```
CREATE TABLE salaries_21 (
    team VARCHAR,
    full_name VARCHAR,
    salary DECIMAL,
    FOREIGN KEY (full_name) REFERENCES active_info(full_name)
);
```