

ETL: NCAA College Football Statistics

Kimberley Reeves, Kevin Clark, Poonam Goel, Rebekah Rowland

Original Project Proposal:

- **Data sources:** We choose to scrape ESPN.com for both full team data and individual player statistics.
- **What Transformations:** Loading website HTML into Jupyter notebooks to be converted into pandas dataframes. Then, converting the data into a JSON and uploading into MongoDB for use as an API later down the road
- **Where is data going to loaded:** Jupyter notebooks, then to MongoDB

Extract:

- We decided to use BeautifulSoup to scrape the URLs that were set up in HTML to get the passing, receiving, and rushing stats for all 130 NCAA teams.
 - **Full Team URL:** <https://www.espn.com/college-football/teams>
 - **Individual Team URLs: 130 Teams**
Example: https://www.espn.com/college-football/team/stats/_/id/2229/florida-international-panthers

Transform:

- The HTML string itself was set up pretty clearly and fairly easy to parse. Using a four loop, we were able to import the individual team URLs into a Pandas Dataframe. From there, we were able to clean up the multiple statistical tables into three clean dataframes: passing, rushing, and receiving. We then imported those data frames into a dictionary, converting the dataframes into JSON strings that were ready to be uploaded.

```
def stats_table_merge(stats_tables):  
    passing=pd.merge(stats_tables[0], stats_tables[1], left_index=True, right_index=True)  
    rushing=pd.merge(stats_tables[2], stats_tables[3], left_index=True, right_index=True)  
    receiving=pd.merge(stats_tables[4], stats_tables[5], left_index=True, right_index=True)  
    team_data={ "passing":passing.to_json(),  
                "rushing":rushing.to_json(),  
                "receiving":receiving.to_json()  
            }  
    return(team_data)
```

Load:

We chose to load our information into MongoDB because we wanted to avoid creating schemas and we had uniform tables which made it easy to upload without having to worry about primary keys.

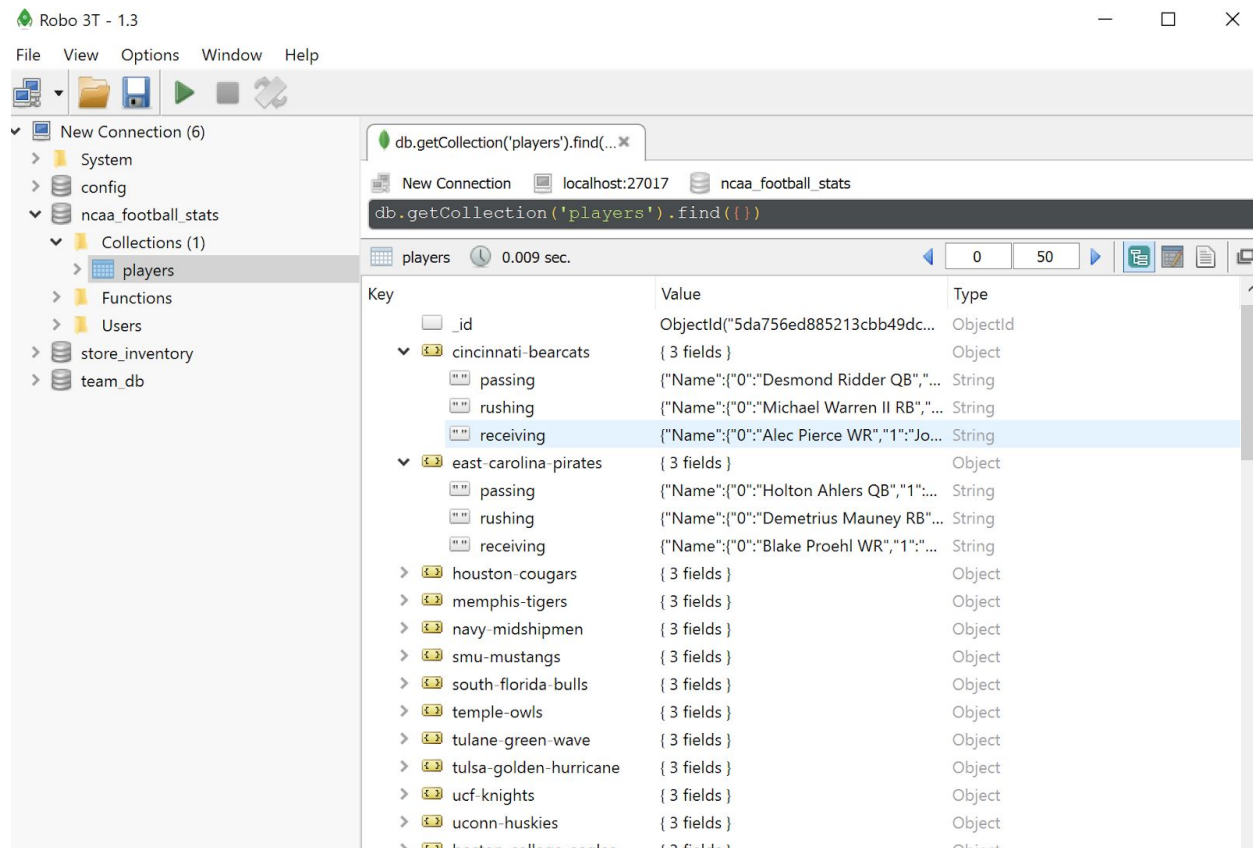
```
client = MongoClient('localhost', 27017)
db = client['passing_db']
collection_passing = db['passing']

with open('C:\\Users\\Poona\\Desktop\\Passing_DataFrame.json') as f:
    file_data = json.load(f)

# use collection_passing.insert(file_data) if pymongo version < 3.0
collection_passing.insert_many(file_data)
client.close()
```

Final MongoDB:

The final collection has the passing, rushing, and receiving statistics saved as JSON strings for 130 NCAA teams.



The screenshot shows the Robo 3T 1.3 application interface. On the left, a tree view shows the database structure: 'ncaa_football_stats' > 'collections (1)' > 'players'. The main window displays the 'players' collection with a query 'db.getCollection('players').find({})' and a result set. The result set is a table with columns 'Key', 'Value', and 'Type'. It lists 130 NCAA teams, each with a '_id' (ObjectId) and three fields: 'passing', 'rushing', and 'receiving' (all String types). The 'receiving' field for the first team, 'cincinnati-bearcats', is highlighted.

Key	Value	Type
_id	ObjectId("5da756ed885213cbb49dc...")	ObjectId
cincinnati-bearcats	{ 3 fields }	Object
passing	{"Name":{"0":"Desmond Ridder QB", "1":...}}	String
rushing	{"Name":{"0":"Michael Warren II RB", "1":...}}	String
receiving	{"Name":{"0":"Alec Pierce WR", "1":"Jo...}}	String
east-carolina-pirates	{ 3 fields }	Object
passing	{"Name":{"0":"Holton Ahlers QB", "1":...}}	String
rushing	{"Name":{"0":"Demetrius Mauney RB", "1":...}}	String
receiving	{"Name":{"0":"Blake Proehl WR", "1":...}}	String
houston-cougars	{ 3 fields }	Object
memphis-tigers	{ 3 fields }	Object
navy-midshipmen	{ 3 fields }	Object
smu-mustangs	{ 3 fields }	Object
south-florida-bulls	{ 3 fields }	Object
temple-owls	{ 3 fields }	Object
tulane-green-wave	{ 3 fields }	Object
tulsa-golden-hurricane	{ 3 fields }	Object
ucf-knights	{ 3 fields }	Object
uconn-huskies	{ 3 fields }	Object
boston-college-eagles	{ 3 fields }	Object