

Search and Filter DataFrames in PySpark HW

Now if's time to put what you've learn into action with a homework assignment!

In case you need it again, here is the link to the documentation for the full list available function in pyspark.sql.functions library: <http://spark.apache.org/docs/latest/api/python/pyspark.sql.html#module-pyspark.sql.functions>

First set up your Spark Session!

Almost so far things first, let's start up our pyspark instance.

```
In [79]: import findspark
findspark.init()
```

```
In [7]: from pyspark.sql import SparkSession
spark=SparkSession.builder.appName("DFAApp").getOrCreate()
```

Read in the DataFrame for this Notebook

We will be continuing to use the fifa19.csv file for this notebook. Make sure that you are writing the correct path to the file.

```
In [12]: df=spark.read.csv("D:\fifa.csv",header=True, inferSchema=True)

df.show()
```

	c_d	ID	Name	Age	Photo	Nationality	Flag	Overall	Potential	Club	...	Composure	Marking	Standing	Tackle
0	158023	L. Messi	31	https://cdn.sofia.org/players/4/19/158023.png	Argentina	https://cdn.sofia.org/flags/92.png	94	94	FC Barcelona	https://cdn.sofia.org/...	13.8	5M	68K	2280	
1	1	20801	Cristiano Ronaldo	33	https://cdn.sofia.org/players/4/19/20801.png	Portugal	https://cdn.sofia.org/flags/98.png	94	94	Juventus	...	95.0	28.0	31.0	180
2	2	190871	Neymar Jr	26	https://cdn.sofia.org/players/4/19/190871.png	Brazil	https://cdn.sofia.org/flags/54.png	92	93	Pars Saint-Germain	...	94.0	27.0	24.0	240
3	3	193080	De Gea	27	https://cdn.sofia.org/players/4/19/193080.png	Spain	https://cdn.sofia.org/flags/45.png	91	93	Manchester United	...	68.0	15.0	21.0	210
4	4	192985	K. De Bruyne	27	https://cdn.sofia.org/players/4/19/192985.png	Belgium	https://cdn.sofia.org/flags/7.png	91	92	Manchester City	...	88.0	68.0	58.0	580
...
18202	18202	238813	J. Lundstram	19	https://cdn.sofia.org/players/4/19/238813.png	England	https://cdn.sofia.org/flags/14.png	47	65	Crewe Alexandra	...	45.0	40.0	48.0	480
18203	18203	243165	N. Christoffersen	19	https://cdn.sofia.org/players/4/19/243165.png	Sweden	https://cdn.sofia.org/flags/46.png	47	63	Trelleborgs FF	...	42.0	22.0	15.0	150
18204	18204	241638	B. Worman	16	https://cdn.sofia.org/players/4/19/241638.png	England	https://cdn.sofia.org/flags/14.png	47	67	Cambridge United	...	41.0	32.0	13.0	130
18205	18205	246268	D. Walker-Rice	17	https://cdn.sofia.org/players/4/19/246268.png	England	https://cdn.sofia.org/flags/14.png	47	66	Tranmere Rovers	...	46.0	20.0	25.0	250
18206	18206	246269	G. Nugent	16	https://cdn.sofia.org/players/4/19/246269.png	England	https://cdn.sofia.org/flags/14.png	46	66	Tranmere Rovers	...	43.0	40.0	43.0	430
18207 rows × 89 columns															

```
In [13]: df.printSchema()
```

```
root
 |-- c_d: integer (nullable = true)
 |-- ID: integer (nullable = true)
 |-- Name: string (nullable = true)
 |-- Age: integer (nullable = true)
 |-- Photo: string (nullable = true)
 |-- Nationality: string (nullable = true)
 |-- Flag: string (nullable = true)
 |-- Overall: integer (nullable = true)
 |-- Potential: integer (nullable = true)
 |-- Club: string (nullable = true)
 |-- Club Logo: string (nullable = true)
 |-- Value: string (nullable = true)
 |-- Wage: string (nullable = true)
 |-- Special: integer (nullable = true)
 |-- Preferred Foot: string (nullable = true)
 |-- International Reputation: integer (nullable = true)
 |-- Weak Foot: integer (nullable = true)
 |-- Skill Moves: integer (nullable = true)
 |-- Work Rate: string (nullable = true)
 |-- Body Type: string (nullable = true)
 |-- Real Face: string (nullable = true)
 |-- Position: string (nullable = true)
 |-- Jersey Number: integer (nullable = true)
 |-- Joined: string (nullable = true)
 |-- Loaned From: string (nullable = true)
 |-- Contract Valid Until: string (nullable = true)
 |-- Height: string (nullable = true)
 |-- Weight: string (nullable = true)
 |-- LS: string (nullable = true)
 |-- RS: string (nullable = true)
 |-- LW: string (nullable = true)
 |-- LF: string (nullable = true)
 |-- CF: string (nullable = true)
 |-- RF: string (nullable = true)
 |-- RW: string (nullable = true)
 |-- CAM: string (nullable = true)
 |-- CAM: string (nullable = true)
 |-- RM: string (nullable = true)
 |-- LM: string (nullable = true)
 |-- LCM: string (nullable = true)
 |-- CM: string (nullable = true)
 |-- RCM: string (nullable = true)
 |-- RM: string (nullable = true)
 |-- LWB: string (nullable = true)
 |-- LW: string (nullable = true)
 |-- CD: string (nullable = true)
 |-- RDM: string (nullable = true)
 |-- RWB: string (nullable = true)
 |-- LB: string (nullable = true)
 |-- LCB: string (nullable = true)
 |-- CB: string (nullable = true)
 |-- RCB: string (nullable = true)
 |-- RB: string (nullable = true)
 |-- Crossing: integer (nullable = true)
 |-- Finishing: integer (nullable = true)
 |-- HeadingAccuracy: integer (nullable = true)
 |-- ShortPassing: integer (nullable = true)
 |-- Volleys: integer (nullable = true)
 |-- Dribbling: integer (nullable = true)
 |-- Curve: integer (nullable = true)
 |-- FKAccuracy: integer (nullable = true)
 |-- LongPassing: integer (nullable = true)
 |-- BallControl: integer (nullable = true)
 |-- Acceleration: integer (nullable = true)
 |-- SprintSpeed: integer (nullable = true)
 |-- Agility: integer (nullable = true)
 |-- Reactions: integer (nullable = true)
 |-- Balance: integer (nullable = true)
 |-- ShotPower: integer (nullable = true)
 |-- Jumping: integer (nullable = true)
 |-- Stamina: integer (nullable = true)
 |-- Strength: integer (nullable = true)
 |-- LongShots: integer (nullable = true)
 |-- Aggression: integer (nullable = true)
 |-- Interceptions: integer (nullable = true)
 |-- Positioning: integer (nullable = true)
 |-- Vision: integer (nullable = true)
 |-- Penalties: integer (nullable = true)
 |-- Composure: integer (nullable = true)
 |-- Marking: integer (nullable = true)
 |-- StandingTackle: integer (nullable = true)
 |-- SlidingTackle: integer (nullable = true)
 |-- GKDividing: integer (nullable = true)
 |-- GKHandling: integer (nullable = true)
 |-- GKkicking: integer (nullable = true)
 |-- GKReflexes: integer (nullable = true)
 |-- Release Clause: string (nullable = true)
```

About this dataframe

The **fifa19.csv** dataset includes a list of all the FIFA 2019 players and their attributes listed below:

- **General:** Age, Nationality, Overall, Potential, Club
- **Metrics:** Value, Wage
- **Player Descriptive:** Preferred Foot, International Reputation, Weak Foot, Skill Moves, Work Rate, Position, Jersey Number, Joined, Loaned From, Contract Valid Until, Height, Weight
- **Position:** LS, ST, RS, LW, LF, CF, RF, RW, LAM, CAM, RAM, LM, LCM, CM, RCM, RM, LWB, LDM, CDM, RDM, RWB, LB, LCB, CB, RCB, RB
- **Other:** Crossing, Finishing, Heading, Accuracy, ShortPassing, Volleys, Dribbling, Curve, FKAccuracy, LongPassing, LongPassing, BallControl, Acceleration, SprintSpeed, Agility, Reactions, Balance, ShotPower, Jumping, Stamina, Strength, LongShots, Aggression, Interceptions, Positioning, Vision, Penalties, Composure, Marking, Standing Tackle, Sliding Tackle, GKDividing, GKHandling, GKkicking, GKPositioning, GKReflexes, and Release Clause.

Source: <https://www.kaggle.com/karangadiya/fifa19>

Use the .toPandas() method to view the first few lines of the dataset so we know what we are working with.

```
In [14]: df.toPandas()
```

```
Out[14]:
```

	c_d	ID	Name	Age	Photo	Nationality	Flag	Overall	Potential	Club	...	Composure	Marking	Standing	Tackle
0	0	158023	L. Messi	31	https://cdn.sofia.org/players/4/19/158023.png	Argentina	https://cdn.sofia.org/flags/92.png	94	94	FC Barcelona	...	96.0	33.0	28.0	280
1	1	20801	Cristiano Ronaldo	33	https://cdn.sofia.org/players/4/19/20801.png	Portugal	https://cdn.sofia.org/flags/98.png	94	94	Juventus	...	95.0	28.0	31.0	180
2	2	190871	Neymar Jr	26	https://cdn.sofia.org/players/4/19/190871.png	Brazil	https://cdn.sofia.org/flags/54.png	92	93	Pars Saint-Germain	...	94.0	27.0	24.0	240
3	3	193080	De Gea	27	https://cdn.sofia.org/players/4/19/193080.png	Spain	https://cdn.sofia.org/flags/45.png	91	93	Manchester United	...	68.0	15.0	21.0	210
4	4	192985	K. De Bruyne	27	https://cdn.sofia.org/players/4/19/192985.png	Belgium	https://cdn.sofia.org/flags/7.png	91	92	Manchester City	...	88.0	68.0	58.0	580
...
18202	18202	238813	J. Lundstram	19	https://cdn.sofia.org/players/4/19/238813.png	England	https://cdn.sofia.org/flags/14.png	47	65	Crewe Alexandra	...	45.0	40.0	48.0	480
18203	18203	243165	N. Christoffersen	19	https://cdn.sofia.org/players/4/19/243165.png	Sweden	https://cdn.sofia.org/flags/46.png	47	63	Trelleborgs FF	...	42.0	22.0	15.0	150
18204	18204	241638	B. Worman	16	https://cdn.sofia.org/players/4/19/241638.png	England	https://cdn.sofia.org/flags/14.png	47	67	Cambridge United	...	41.0	32.0	13.0	130
18205	18205	246268	D. Walker-Rice	17	https://cdn.sofia.org/players/4/19/246268.png	England	https://cdn.sofia.org/flags/14.png	47	66	Tranmere Rovers	...	46.0	20.0	25.0	250
18206	18206	246269	G. Nugent	16	https://cdn.sofia.org/players/4/19/246269.png	England	https://cdn.sofia.org/flags/14.png	46	66	Tranmere Rovers	...	43.0	40.0	43.0	430
18207 rows × 89 columns															

Now print the schema of the dataset so we can see the data types of all the variables.

```
In [13]: df.printSchema()
```

```
root
 |-- c_d: integer (nullable = true)
 |-- ID: integer (nullable = true)
 |-- Name: string (nullable = true)
 |-- Age: integer (nullable = true)
 |-- Photo: string (nullable = true)
 |-- Nationality: string (nullable = true)
 |-- Flag: string (nullable = true)
 |-- Overall: integer (nullable = true)
 |-- Potential: integer (nullable = true)
 |-- Club: string (nullable = true)
 |-- Club Logo: string (nullable = true)
 |-- Value: string (nullable = true)
 |-- Wage: string (nullable = true)
 |-- Special: integer (nullable = true)
 |-- Preferred Foot: string (nullable = true)
 |-- International Reputation: integer (nullable = true)
 |-- Weak Foot: integer (nullable = true)
 |-- Skill Moves: integer (nullable = true)
 |-- Work Rate: string (nullable = true)
 |-- Body Type: string (nullable = true)
 |-- Real Face: string (nullable = true)
 |-- Position: string (nullable = true)
 |-- Jersey Number: integer (nullable = true)
 |-- Joined: string (nullable = true)
 |-- Loaned From: string (nullable = true)
 |-- Contract Valid Until: string (nullable = true)
 |-- Height: string (nullable = true)
 |-- Weight: string (nullable = true)
 |-- LS: string (nullable = true)
 |-- RS: string (nullable = true)
 |-- LW: string (nullable = true)
 |-- LF: string (nullable = true)
 |-- CF: string (nullable = true)
 |-- RF: string (nullable = true)
 |-- RW: string (nullable = true)
 |-- CAM: string (nullable = true)
 |-- CAM: string (nullable = true)
 |-- RM: string (nullable = true)
 |-- LM: string (nullable = true)
 |-- LCM: string (nullable = true)
 |-- CM: string (nullable = true)
 |-- RCM: string (nullable = true)
 |-- RM: string (nullable = true)
 |-- LWB: string (nullable = true)
 |-- LW: string (nullable = true)
 |-- CD: string (nullable = true)
 |-- RDM: string (nullable = true)
 |-- RWB: string (nullable = true)
 |-- LB: string (nullable = true)
 |-- LCB: string (nullable = true)
 |-- CB: string (nullable = true)
 |-- RCB: string (nullable = true)
 |-- RB: string (nullable = true)
 |-- Crossing: integer (nullable = true)
 |-- Finishing: integer (nullable = true)
 |-- HeadingAccuracy: integer (nullable = true)
 |-- ShortPassing: integer (nullable = true)
 |-- Volleys: integer (nullable = true)
 |-- Dribbling: integer (nullable = true)
 |-- Curve: integer (nullable = true)
 |-- FKAccuracy: integer (nullable = true)
 |-- LongPassing: integer (nullable = true)
 |-- BallControl: integer (nullable = true)
 |-- Acceleration: integer (nullable = true)
 |-- SprintSpeed: integer (nullable = true)
 |-- Agility: integer (nullable = true)
 |-- Reactions: integer (nullable = true)
 |-- Balance: integer (nullable = true)
 |-- ShotPower: integer (nullable = true)
 |-- Jumping: integer (nullable = true)
 |-- Stamina: integer (nullable = true)
 |-- Strength: integer (nullable = true)
 |-- LongShots: integer (nullable = true)
 |-- Aggression: integer (nullable = true)
 |-- Interceptions: integer (nullable = true)
 |-- Positioning: integer (nullable = true)
 |-- Vision: integer (nullable = true)
 |-- Penalties: integer (nullable = true)
 |-- Composure: integer (nullable = true)
 |-- Marking: integer (nullable = true)
 |-- StandingTackle: integer (nullable = true)
 |-- SlidingTackle: integer (nullable = true)
 |-- GKDividing: integer (nullable = true)
 |-- GKHandling: integer (nullable = true)
 |-- GKkicking: integer (nullable = true)
 |-- GKReflexes: integer (nullable = true)
 |-- Release Clause: string (nullable = true)
```

Now let's get started!

First things first..... import the pyspark sql functions library

Since we know we will be using it a lot.

```
In [41]: import pyspark.sql.functions as f
```

1. Select the Name and Position of each player in the dataframe

```
In [57]: df.createOrReplaceTempView("tempview")
res=spark.sql("select Name,Position from tempview")
res.toPandas()
```

```
Out[57]:
```

	Name	Position
0	L. Messi	RF
1	Cristiano Ronaldo	ST
2	Neymar Jr	LW
3	De Gea	GM
4	K. De Bruyne	RCM
...
18202	J. Lundstram	CM
18203	N. Christoffersen	ST
18204	B. Worman	ST
18205	D. Walker-Rice	RW
18206	G. Nugent	CM
18207 rows × 2 columns		

1.1 Display the same results from above sorted by the players names

```
In [71]: res=spark.sql("select Name,Position from tempview order by Name")
res.toPandas()
```

```
Out[71]:
```

	Name	Position
0	A. Abiang	ST
1	A. Abdelnour	LB
2	A. Abdelmonem	CB
3	A. Abdi	CM
4	A. Abdul Jabbar	ST
...
18202	Edgar Plano	LM
18203	Edgar Valentin	CDM
18204	Edgar Valentin	CDM
18205	Edgar Valentin	CDM
18206	Edgar Valentin	LCB
18207 rows × 2 columns		

2. Select only the players who belong to a club beginning with FC

```
In [26]: res=spark.sql("select Name,Club from tempview where club like 'FC%'")
res.toPandas()
```

```
Out[26]:
```

	Name	Club
0	L. Messi	FC Barcelona
1	L. Suñez	FC Barcelona
2	R. Lewandowski	FC Bayern Mfchen
3	M. ter Stegen	FC Barcelona
4	Sergio Busquets	FC Barcelona
...
1002	M. Hamper	FC Wfzburger Kickers
1003	O. Dzonjagic	FC Thun
1004	O. Olsen	FC Midtjylland
1005	N. Stephan	FC Wfzburger Kickers
1006	M. Finne Wfz	FC Nordsjlland
1007 rows × 2 columns		

3. Who is the oldest player in the dataset and how old are they?

Display only the name and age of the oldest player.

```
In [22]: res=spark.sql("select ID,Name,Age from tempview order by Age desc")
res.show(1)
```

```
+-----+
|ID| Name|Age|
+-----+
|148029|O. Pföz| 45|
+-----+
```

only showing top 1 row

4. Select only the following players from the dataframe:

- L. Messi
- Cristiano Ronaldo

```
In [63]: res=spark.sql("select Name from tempview where Name like '%L. Messi' or Name like '%Cristiano Ronaldo%'")
res.toPandas()
```

```
Out[63]:
```

	Name
0	L. Messi
1	Cristiano Ronaldo

5. Can you select the first character from the Release Clause variable which indicates the currency used?

```
In [68]: df.select("Release Clause",df["Release Clause"].substr(1,1)).toPandas()
```

```
Out[68]:
```

	Release Clause	substring(Release Clause, 1, 1)
0	€26.5M	€
1	€27.1M	€
2	€28.1M	€
3	€38.6M	€
4	€36.4M	€
...
18202	€3K	€
18203	€1K	€
18204	€65K	€
18205	€3K	€
18206	€65K	€
18207 rows × 2 columns		

6. Can you select only the players who are over the age of 40?

```
In [35]: res=spark.sql("select Name,Age from tempview where Age>40")
res.toPandas()
```

```
Out[35]:
```

	Name	Age
0	J. Villar	41
1	B. Nivet	41
2	O. Pfoz	45
3	C. Muek	41
4	S. Nazranzi	42
5	H. Malsami	41
6	M. Tyler	41
7	T. Warner	44
8	K. Pilkington	44

That's for now... Great Job!