

```
In [1]: ▶ 1 import pybaseball as pyb
2 from pybaseball import statcast, pitching_stats, playerid_lookup, stat
3 import numpy as np
4 import math
5 import pandas as pd
6 import seaborn as sns
7 import matplotlib.pyplot as plt
8 %matplotlib inline
9 import glob
10 import os
11 import re
12 import unicodedata
13 from datetime import datetime
14 from itertools import groupby
15 from operator import itemgetter
16 from fuzzywuzzy import process
```

C:\Users\johns\anaconda3\Lib\site-packages\pandas\core\arrays\masked.py:6
0: UserWarning: Pandas requires version '1.3.6' or newer of 'bottleneck'
(version '1.3.5' currently installed).
from pandas.core import (

In [2]:

▶

```
1 hist_tj_df = pd.read_csv('hist_tj.csv', index_col=0)
2 hist_tj_df
```

Out[2]:

	Name	Age	Year	Throws	IP	G	GS	CG	SHO	sDR	Career Start	Career End	Inact Yea
240	steve blass	31	1973	1	88.2	23	18	1	0	1	1972	1973	
285	eddie fisher	36	1973	1	117.2	32	16	2	0	9	1972	1973	
385	milt pappas	34	1973	1	162.0	30	29	1	1	8	1972	1973	
463	steve arlin	28	1974	1	107.2	27	22	2	0	9	1972	1974	
589	ernie mcanally	27	1974	1	128.2	25	21	5	2	6	1972	1974	
...	
15334	zack wheeler	33	2023	1	192.0	32	32	0	0	0	2013	2023	[20 20
15339	trevor williams	31	2023	1	144.1	30	30	0	0	0	2016	2023	
15346	alex wood	32	2023	0	97.2	29	12	0	0	1	2013	2023	
15348	brandon woodruff	30	2023	1	67.0	11	11	1	1	0	2017	2023	
15350	ryan yarbrough	31	2023	0	89.2	25	9	0	0	0	2018	2023	

1244 rows × 20 columns



In [3]: 1 hist_tj_df.info()

```
<class 'pandas.core.frame.DataFrame'>
Index: 1244 entries, 240 to 15350
Data columns (total 20 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Name                   1244 non-null   object
1   Age                    1244 non-null   int64
2   Year                   1244 non-null   int64
3   Throws                 1244 non-null   int64
4   IP                     1244 non-null   float64
5   G                      1244 non-null   int64
6   GS                     1244 non-null   int64
7   CG                     1244 non-null   int64
8   SHO                    1244 non-null   int64
9   SDR                    1244 non-null   int64
10  Career Start           1244 non-null   int64
11  Career End             1244 non-null   int64
12  Inactive Years         1244 non-null   object
13  Surgery                 1244 non-null   float64
14  TJ Surgery Date        1244 non-null   object
15  Surgeon(s)             197 non-null    object
16  Country                 268 non-null    object
17  Level                   268 non-null    object
18  Total_IP                1244 non-null   float64
19  TJ Surgery Year         1244 non-null   object
dtypes: float64(3), int64(10), object(7)
memory usage: 204.1+ KB
```

```
In [4]: 1 pitcher_data_df = hist_tj_df
        2 pitcher_data_df
```

Out[4]:

	Name	Age	Year	Throws	IP	G	GS	CG	SHO	sDR	Career Start	Career End	Inact Year
240	steve blass	31	1973	1	88.2	23	18	1	0	1	1972	1973	
285	eddie fisher	36	1973	1	117.2	32	16	2	0	9	1972	1973	
385	milt pappas	34	1973	1	162.0	30	29	1	1	8	1972	1973	
463	steve arlin	28	1974	1	107.2	27	22	2	0	9	1972	1974	
589	ernie mcanally	27	1974	1	128.2	25	21	5	2	6	1972	1974	
...	
15334	zack wheeler	33	2023	1	192.0	32	32	0	0	0	2013	2023	[2020
15339	trevor williams	31	2023	1	144.1	30	30	0	0	0	2016	2023	
15346	alex wood	32	2023	0	97.2	29	12	0	0	1	2013	2023	
15348	brandon woodruff	30	2023	1	67.0	11	11	1	1	0	2017	2023	
15350	ryan yarbrough	31	2023	0	89.2	25	9	0	0	0	2018	2023	

1244 rows × 20 columns



Goal: Take previous DF and attach identification ('key_mlbam') back to the players.

```
In [5]: ▶ 1 # Function to fetch key_mlbam
2 def fetch_key_mlbam(row):
3     try:
4         # Splitting the name into first and last name
5         first_name, last_name = row['Name'].split(' ')[0], ' '.join(row['Name'].split(' ')[1:])
6         # Fetching player ID
7         player_id_df = playerid_lookup(last_name, first_name)
8         # Assuming the first result is the correct one, adjust as needed
9         key_mlbam = player_id_df.iloc[0]['key_mlbam']
10        return key_mlbam
11    except Exception as e:
12        print(f"Error fetching key_mlbam for {row['Name']}: {e}")
13        return pd.NA
14
15 # Apply the function to each row and create a new column 'key_mlbam'
16 pitcher_data_df['key_mlbam'] = pitcher_data_df.apply(fetch_key_mlbam, axis=1)
17
18 # Display the updated DataFrame
19 pitcher_data_df.head()
```



```

In [6]: 1 names_list = ['blue moon odom', 'jr richard', 'silvio martinez', 'will
2 'john henry johnson', 'alejandro pena', 'joaquin andujar', 'oil can bo
3 'pascual perez', 'jose deleon', 'jose mesa', 'jose de jesus',
4 'jose guzman', 'jose bautista', 'angel miranda', 'dennis martinez',
5 'francisco cordova', 'juan guzman', 'carlos perez', 'hipolito pichardo
6 'jose rosado', 'jose silva', 'osvaldo fernandez', 'ramon martinez',
7 'jose mercedes', 'vladimir nunez', 'jose rijo', 'ruben quevedo',
8 'salomon torres', 'jesus sanchez', 'wilson alvarez', 'joaquin benoit',
9 'geremi gonzalez', 'sunwoo kim', 'jose lima', 'gustavo chacin',
10 'orlando hernandez', 'byunghyun kim', 'victor santos', 'jae weong seo'
11 'julian tavarez', 'victor zambrano', 'shawn chacon', 'runelvys hernand
12 'hungchih kuo', 'odalis perez', 'jose contreras', 'chan ho park',
13 'horacio ramirez', 'oliver perez', 'ryan rowlandsmith', 'dj carrasco',
14 'livan hernandez', 'rodrigo lopez', 'joel pineiro', 'jojo reyes', 'jav
15 'freddy garcia', 'ramon ortiz', 'jonathan sanchez', 'chienming wang',
16 'erik bedard', 'aj burnett', 'felix doubront', 'wandy rodriguez',
17 'cj wilson', 'jose fernandez', 'roberto hernandez', 'jon niese',
18 'vidal nuno iii', 'alfredo simon', 'henderson alvarez iii', 'ra dickey
19 'aj griffin', 'ubaldo jimenez', 'weiyin chen', 'bartolo colon', 'roeni
20 'jaime garcia', 'miguel gonzalez', 'felix hernandez', 'hector noesi',
21 'edinson volquez', 'ivan nova', 'jose alvarez', 'jhoulys chacin', 'ja
22 'jorge lopez', 'carlos martinez', 'hector santiago', 'reynaldo lopez',
23 'lance mccullers jr', 'anibal sanchez', 'sandy alcantara', 'jaime barr
24 'matthew boyd', 'nestor cortes', 'domingo german', 'carlos hernandez',
25 'jesus luzardo', 'german marquez', 'nick martinez', 'martin perez', 'j
26 'carlos rodon', 'eduardo rodriguez', 'hyun jin ryu', 'jose suarez', 'r
27 'julio teheran', 'jose urena', 'julio urias', 'jose urquidy']
28
29 # Target name you are trying to match
30 all_player_names = ['Blue Moon Odom', 'J. R. Richard', 'Silvio Martíne
31 'Willie Hernández', 'John Henry Johnson', 'Alejand
32 'Joaquín Andújar', 'Oil Can Boyd', 'Pascual Pérez'
33 'José Mesa', 'José DeJesús', 'José Guzmán', 'José
34 'Ángel Miranda', 'Dennis Martínez', 'Francisco Cór
35 'Juan Guzmán', 'Carlos Pérez', 'Hipólito Pichardo'
36 'José Silva', 'Osvaldo Fernández', 'Ramón Martínez
37 'Vladimir Núñez', 'José Rijo', 'Rubén Quevedo', 'S
38 'Jesús Sánchez', 'Wilson Álvarez', 'Joaquín Benoit
39 'Sun-woo Kim', 'José Lima', 'Gustavo Chacín', 'Orl
40 'Byung-hyun Kim', 'Víctor Santos', 'Jae Weong Seo'
41 'Julián Tavárez', 'Víctor Zambrano', 'Shawn Chacón
42 'Runelvys Hernández', 'Hong-Chih Kuo', 'Odalis Pér
43 'José Contreras', 'Chan Ho Park', 'Horacio Ramírez'
44 'Ryan Rowland-Smith', 'D. J. Carrasco', 'Liván Her
45 'Rodrigo López', 'Joel Piñeiro', 'Jo-Jo Reyes', 'J
46 'Freddy García', 'Ramón Ortiz', 'Jonathan Sánchez'
47 'Érik Bédard', 'A. J. Burnett', 'Félix Doubront',
48 'C. J. Wilson', 'José Fernández', 'Roberto Hernánde
49 'Jon Niese', 'Vidal Nuño', 'Alfredo Simón', 'Hende
50 'R. A. Dickey', 'A. J. Griffin', 'Ubaldo Jiménez',
51 'Bartolo Colón', 'Roenis Elías', 'Jaime García', 'I
52 'Félix Hernández', 'Héctor Noesí', 'Edinson Vólque
53 'José Álvarez', 'Jhoulys Chacín', 'J.A. Happ', 'Jon
54 'Carlos Martínez', 'Hector Santiago', 'Reynaldo Ló
55 'Lance McCullers Jr.', 'Aníbal Sánchez', 'Sandy Al
56 'José Berríos', 'Matthew Boyd', 'Nestor Cortes', 'I
57 'Carlos Hernández', 'Pablo López', 'Jesús Luzardo']

```

```
58         'Germán Márquez', 'Martín Pérez', 'Carlos Rodón', '  
59         'Hong-Chih Kuo', 'Edinson Vólquez', 'Jose Alvarez'  
60  
61 # Function to find the best match for each name in names_list  
62 def find_best_matches(names_list, all_player_names):  
63     best_matches = {}  
64     for name in names_list:  
65         match = process.extractOne(name, all_player_names)  
66         best_matches[name] = match  
67     return best_matches  
68  
69 # Get best matches  
70 best_matches = find_best_matches(names_list, all_player_names)  
71  
72 # Print best matches  
73 for name, match in best_matches.items():  
74     print(f"Original: {name}, Best Match: {match[0]}, Score: {match[1]}
```



```
In [7]: 1 updated_names = {  
2     "blue moon odom": "Blue Moon Odom",  
3     "jr richard": "J. R. Richard",  
4     "silvio martinez": "Silvio Martínez",  
5     "willie hernandez": "Willie Hernández",  
6     "john henry johnson": "John Henry Johnson",  
7     "alejandro pena": "Alejandro Peña",  
8     "joaquin andujar": "Joaquín Andújar",  
9     "oil can boyd": "Oil Can Boyd",  
10    "pascual perez": "Pascual Pérez",  
11    "jose deleon": "José DeLeón",  
12    "jose mesa": "José Mesa",  
13    "jose de jesus": "José DeJesús",  
14    "jose guzman": "José Guzmán",  
15    "jose bautista": "José Bautista",  
16    "angel miranda": "Ángel Miranda",  
17    "dennis martinez": "Dennis Martínez",  
18    "francisco cordova": "Francisco Córdova",  
19    "juan guzman": "Juan Guzmán",  
20    "carlos perez": "Carlos Pérez",  
21    "hipolito pichardo": "Hipólito Pichardo",  
22    "jose rosado": "José Rosado",  
23    "jose silva": "José Silva",  
24    "osvaldo fernandez": "Osvaldo Fernández",  
25    "ramon martinez": "Ramón Martínez",  
26    "jose mercedes": "José Mercedes",  
27    "vladimir nunez": "Vladimir Núñez",  
28    "jose rijo": "José Rijo",  
29    "ruben quevedo": "Rubén Quevedo",  
30    "salomon torres": "Salomón Torres",  
31    "jesus sanchez": "Jesús Sánchez",  
32    "wilson alvarez": "Wilson Álvarez",  
33    "joaquin benoit": "Joaquín Benoit",  
34    "geremi gonzalez": "Geremi González",  
35    "sunwoo kim": "Sun-woo Kim",  
36    "jose lima": "José Lima",  
37    "gustavo chacin": "Gustavo Chacín",  
38    "orlando hernandez": "Orlando Hernández",  
39    "byunghyun kim": "Byung-hyun Kim",  
40    "victor santos": "Víctor Santos",  
41    "jae weong seo": "Jae Weong Seo",  
42    "julian tavarez": "Julián Tavárez",  
43    "victor zambrano": "Víctor Zambrano",  
44    "shawn chacon": "Shawn Chacón",  
45    "runelvys hernandez": "Runelvys Hernández",  
46    "hungchih kuo": "Hong-Chih Kuo",  
47    "odalis perez": "Odalis Pérez",  
48    "jose contreras": "José Contreras",  
49    "chan ho park": "Chan Ho Park",  
50    "horacio ramirez": "Horacio Ramírez",  
51    "oliver perez": "Óliver Pérez",  
52    "ryan rowlandsmith": "Ryan Rowland-Smith",  
53    "dj carrasco": "D. J. Carrasco",  
54    "livan hernandez": "Liván Hernández",  
55    "rodrigo lopez": "Rodrigo López",  
56    "joel pineiro": "Joel Piñeiro",  
57    "jojo reyes": "Jo-Jo Reyes",
```

```
58 "javier vazquez": "Javier Vázquez",
59 "freddy garcia": "Freddy García",
60 "ramon ortiz": "Ramón Ortiz",
61 "jonathan sanchez": "Jonathan Sánchez",
62 "chienming wang": "Chien-Ming Wang",
63 "erik bedard": "Érik Bédard",
64 "aj burnett": "A. J. Burnett",
65 "felix doubront": "Félix Doubront",
66 "wandy rodriguez": "Wandy Rodríguez",
67 "cj wilson": "C. J. Wilson",
68 "jose fernandez": "José Fernández",
69 "roberto hernandez": "Roberto Hernández",
70 "jon niese": "Jon Niese",
71 "vidal nuno iii": "Vidal Nuño",
72 "alfredo simon": "Alfredo Simón",
73 "henderson alvarez iii": "Henderson Álvarez",
74 "ra dickey": "R. A. Dickey",
75 "aj griffin": "A. J. Griffin",
76 "ubaldo jimenez": "Ubaldo Jiménez",
77 "weiyin chen": "Wei-Yin Chen",
78 "bartolo colon": "Bartolo Colón",
79 "roenis elias": "Roenis Elías",
80 "jaime garcia": "Jaime García",
81 "miguel gonzalez": "Miguel González",
82 "felix hernandez": "Félix Hernández",
83 "hector noesi": "Héctor Noesí",
84 "edinson volquez": "Edinson Vólquez",
85 "ivan nova": "Iván Nova",
86 "jose alvarez": "Jose Alvarez",
87 "jhoulys chacin": "Jhoulys Chacín",
88 "ja happ": "J.A. Happ",
89 "jorge lopez": "Jorge López",
90 "carlos martinez": "Carlos Martínez",
91 "hector santiago": "Hector Santiago",
92 "reynaldo lopez": "Reynaldo López",
93 "lance mccullers jr": "Lance McCullers Jr.",
94 "anibal sanchez": "Aníbal Sánchez",
95 "sandy alcantara": "Sandy Alcántara",
96 "jaime barria": "Jaime Barría",
97 "jose berrios": "José Berríos",
98 "matthew boyd": "Matthew Boyd",
99 "nestor cortes": "Nestor Cortes",
100 "domingo german": "Domingo Germán",
101 "carlos hernandez": "Carlos Hernández",
102 "pablo lopez": "Pablo López",
103 "jesus luzardo": "Jesús Luzardo",
104 "german marquez": "Germán Márquez",
105 "nick martinez": "Dennis Martínez",
106 "martin perez": "Martín Pérez",
107 "jose quintana": "José Quintana",
108 "carlos rodon": "Carlos Rodón",
109 "eduardo rodriguez": "Eduardo Rodríguez",
110 "hyun jin ryu": "Hyun-jin Ryu",
111 "jose suarez": "José Suarez",
112 "ranger suarez": "Ranger Suárez",
113 "julio teheran": "Julio Teheran",
114 "jose urena": "José Ureña",
```

```

115     "julio urias": "Julio Urías",
116     "jose urquidy": "José Urquidy",
117 }
118
119 pitcher_data_df['Name'] = pitcher_data_df['Name'].map(updated_names).f
120 pitcher_data_df

```

```

In [8]: ▶ 1 # Function to fetch key_mlbam
2 def fetch_key_mlbam(row):
3     try:
4         # Splitting the name into first and last name
5         first_name, last_name = row['Name'].split(' ')[0], ' '.join(ro
6         # Fetching player ID
7         player_id_df = playerid_lookup(last_name, first_name)
8         # Assuming the first result is the correct one, adjust as neces
9         key_mlbam = player_id_df.iloc[0]['key_mlbam']
10        return key_mlbam
11    except Exception as e:
12        print(f"Error fetching key_mlbam for {row['Name']}: {e}")
13        return pd.NA
14
15    # Apply the function to each row and create a new column 'key_mlbam'
16    pitcher_data_df['key_mlbam'] = pitcher_data_df.apply(fetch_key_mlbam,
17
18    # Display the updated DataFrame
19    pitcher_data_df.head()

```

From here, looked up player names, referenced back to baseball-reference.com to confirm the player was correct, and manually entered player keys.

```

In [9]: ▶ 1 playerid_lookup('teheran', fuzzy=True)

```

```

In [10]: ▶ 1 teheran_key = '527054'
2 pitcher_data_df.loc[pitcher_data_df['Name'] == "Julio Teheran", 'key_m

```

```

In [11]: ▶ 1 odom_key = '119935'
2 pitcher_data_df.loc[pitcher_data_df['Name'] == "Blue Moon Odom", 'key_

```

```

In [12]: ▶ 1 richard_key = '121145'
2 pitcher_data_df.loc[pitcher_data_df['Name'] == "J. R. Richard", 'key_m

```

```

In [13]: ▶ 1 oil_key = '111312'
2 pitcher_data_df.loc[pitcher_data_df['Name'] == "Oil Can Boyd", 'key_m

```

```
In [14]: 1 seo_key = '150242'
        2 pitcher_data_df.loc[pitcher_data_df['Name'] == "Jae Weong Seo", 'key_m

In [15]: 1 kuo_key = '425539'
        2 pitcher_data_df.loc[pitcher_data_df['Name'] == "Hong-Chih Kuo", 'key_m

In [16]: 1 park_key = '120221'
        2 pitcher_data_df.loc[pitcher_data_df['Name'] == "Chan Ho Park", 'key_ml

In [17]: 1 djcarrasco_key = '425647'
        2 pitcher_data_df.loc[pitcher_data_df['Name'] == "D. J. Carrasco", 'key_

In [18]: 1 burnett_key = '150359'
        2 pitcher_data_df.loc[pitcher_data_df['Name'] == "A. J. Burnett", 'key_m

In [19]: 1 wilson_key = '450351'
        2 pitcher_data_df.loc[pitcher_data_df['Name'] == "C. J. Wilson", 'key_ml

In [20]: 1 niese_key = '477003'
        2 pitcher_data_df.loc[pitcher_data_df['Name'] == "Jon Niese", 'key_mlbam

In [21]: 1 dickey_key = '285079'
        2 pitcher_data_df.loc[pitcher_data_df['Name'] == "R. A. Dickey", 'key_ml

In [22]: 1 griffin_key = '456167'
        2 pitcher_data_df.loc[pitcher_data_df['Name'] == "A. J. Griffin", 'key_m

In [23]: 1 alvarez_key = '571439'
        2 pitcher_data_df.loc[pitcher_data_df['Name'] == "Jose Alvarez", 'key_ml

In [24]: 1 happ_key = '457918'
        2 pitcher_data_df.loc[pitcher_data_df['Name'] == "J.A. Happ", 'key_mlbam

In [25]: 1 santiago_key = '502327'
        2 pitcher_data_df.loc[pitcher_data_df['Name'] == "Hector Santiago", 'key

In [26]: 1 mccullers_key = '621121'
        2 pitcher_data_df.loc[pitcher_data_df['Name'] == "Lance McCullers Jr.",
        <img alt="Horizontal scrollbar" data-bbox="261 758 918 771"/>

In [27]: 1 mattboyd_key = '571510'
        2 pitcher_data_df.loc[pitcher_data_df['Name'] == "Matthew Boyd", 'key_ml
```

Nestor Cortes, José DeJesús

Can't find. Drop.

```
In [28]: 1 pitcher_data_df = pitcher_data_df[pitcher_data_df['Name'] != 'Nestor C...
```

```
In [29]: 1 pitcher_data_df = pitcher_data_df[pitcher_data_df['Name'] != 'José DeJ...
```

```
In [30]: 1 pitcher_data_df['key_mlbam'].isna().value_counts()
```

```
Out[30]: key_mlbam
False    1239
True       3
Name: count, dtype: int64
```

(Within the saved CSV, here are no True value counts. All ID's are accounted for.)

```
In [31]: 1 pitcher_data_df.info()
```

...

```
In [32]: 1 """
2 pitcher_data_df.to_csv('pitcher_key_df.csv')
3 """
```

```
Out[32]: "\npitcher_data_df.to_csv('pitcher_key_df.csv')\n"
```

```
In [33]: 1 playerid_lookup('ohtani', 'shohei')
```

```
Out[33]:
```

	name_last	name_first	key_mlbam	key_retro	key_bbreff	key_fangraphs	mlb_played_first
0	ohtani	shohei	660271	ohtas001	ohtansh01	19755	2018.0

```
In [34]: 1 ohtani_stats = statcast_pitcher('2018-01-01', '2023-10-01', 660271)
2 ohtani_stats
```

...

```
In [35]: 1 ohtani_stats.info()
```

...

```
In [36]: 1 ohtani_stats['game_date'].value_counts()
```

...

```
In [37]: 1 ohtani_stats['game_date'].unique
```

...

```
In [38]: 1 ohtani_start = ohtani_stats[ohtani_stats['game_date'] == '2021-09-03']
2 ohtani_start
```

...

```
In [39]: 1 ohtani_start.columns
```

...

In [40]: 1 ohtani_start.info()

...

This kind of info is great!

Will filter down columns to information only applicable to testing for TJ surgery.

Need to filter years to only include regular season schedule.

This took some manual input.

```
In [41]: 1 """
2 season_dates
3 2008: ('2008-03-25', '2008-09-30'),
4 2009: ('2009-04-05', '2009-10-06'),
5 2010: ('2010-04-04', '2010-10-03'),
6 2011: ('2011-03-31', '2011-09-28'),
7 2012: ('2012-03-28', '2012-10-03'),
8 2013: ('2013-03-31', '2013-09-30'),
9 2014: ('2014-03-22', '2014-09-28'),
10 2015: ('2015-04-05', '2015-10-04'),
11 2016: ('2016-04-03', '2016-10-02'),
12 2017: ('2017-04-02', '2017-10-01'),
13 2018: ('2018-03-29', '2018-10-01'),
14 2019: ('2019-03-20', '2019-09-29'),
15 2020: ('2020-07-23', '2020-09-27'),
16 2021: ('2021-04-01', '2021-10-03'),
17 2022: ('2022-04-07', '2022-10-05'),
18 2023: ('2023-03-30', '2023-10-01')
19 """
```

```
Out[41]: "\nseason_dates\n    2008: ('2008-03-25', '2008-09-30'),\n    2009: ('2009-04-05', '2009-10-06'),\n    2010: ('2010-04-04', '2010-10-03'),\n    2011: ('2011-03-31', '2011-09-28'),\n    2012: ('2012-03-28', '2012-10-03'),\n    2013: ('2013-03-31', '2013-09-30'),\n    2014: ('2014-03-22', '2014-09-28'),\n    2015: ('2015-04-05', '2015-10-04'),\n    2016: ('2016-04-03', '2016-10-02'),\n    2017: ('2017-04-02', '2017-10-01'),\n    2018: ('2018-03-29', '2018-10-01'),\n    2019: ('2019-03-20', '2019-09-29'),\n    2020: ('2020-07-23', '2020-09-27'),\n    2021: ('2021-04-01', '2021-10-03'),\n    2022: ('2022-04-07', '2022-10-05'),\n    2023: ('2023-03-30', '2023-10-01')\n"
```

```

In [42]: ▶ 1 """
2 all_player_stats = []
3
4 # Loop through each MLBAM ID in the pitcher_data_df
5 for key_mlbam in pitcher_data_df['key_mlbam']:
6     # Fetch pitching stats from statcast
7     player_stats = statcast_pitcher('2018-03-29', '2018-10-01', key_mlbam)
8     # Append the fetched stats to the list
9     all_player_stats.append(player_stats)
10
11 # Concatenate all DataFrames into a single DataFrame
12 all_2018_stats_df = pd.concat(all_player_stats, ignore_index=True)
13
14 all_2018_stats_df.head()
15 """

```

```

Out[42]: "\nall_player_stats = []\n\n# Loop through each MLBAM ID in the pitcher_data_df\nfor key_mlbam in pitcher_data_df['key_mlbam']:\n    # Fetch pitching stats from statcast\n    player_stats = statcast_pitcher('2018-03-29', '2018-10-01', key_mlbam)\n    # Append the fetched stats to the list\n    all_player_stats.append(player_stats)\n\n# Concatenate all DataFrames into a single DataFrame\nall_2018_stats_df = pd.concat(all_player_stats, ignore_index=True)\n\nall_2018_stats_df.head()\n"

```

Start grouping data by game date and pitcher.

Later, further condense data to season and pitcher, while retaining important information about the number of different types of pitches, the averages of those pitches velocity, etc.

Was not done here due to sheer size of each season's file.


```

In [44]: 1 """
2 # Group by 'game_date' and 'pitcher' to calculate the total pitches
3 total_pitches = all_2018_stats_df.groupby(['game_date', 'pitcher']).size
4
5 # Group by 'game_date', 'pitcher', and 'pitch_type' to calculate the sum
6 total_pitches_by_type = all_2018_stats_df.groupby(['game_date', 'pitcher', 'pitch_type']).size
7
8 # Calculate averages of the specified metrics for each pitch type, grouped by game_date and pitcher
9 avg_metrics = all_2018_stats_df.groupby(['game_date', 'pitcher', 'pitch_type']).agg({
10     'release_speed': 'mean',
11     'release_pos_x': 'mean',
12     'release_pos_z': 'mean',
13     'spin_dir': 'mean',
14     'vx0': 'mean',
15     'vy0': 'mean',
16     'vz0': 'mean',
17     'ax': 'mean',
18     'ay': 'mean',
19     'az': 'mean',
20     'effective_speed': 'mean',
21     'release_spin_rate': 'mean',
22     'release_extension': 'mean',
23     'release_pos_y': 'mean',
24     'spin_axis': 'mean'
25 }).reset_index()
26
27 # Merging total pitches and total pitches by type back
28 grouped_2018_df = total_pitches.merge(total_pitches_by_type, on=['game_date', 'pitcher'], how='left')
29 grouped_2018_df = final_df.merge(avg_metrics, on=['game_date', 'pitcher', 'pitch_type'], how='left')
30 """

```

```

Out[44]: "\n# Group by 'game_date' and 'pitcher' to calculate the total pitches\n
total_pitches = all_2018_stats_df.groupby(['game_date', 'pitcher']).size
().reset_index(name='total_pitches')\n\n# Group by 'game_date', 'pitcher', and 'pitch_type' to calculate the sum total of each pitch_type\n
total_pitches_by_type = all_2018_stats_df.groupby(['game_date', 'pitcher', 'pitch_type']).size().reset_index(name='count_by_pitch_type')\n\n# Calculate averages of the specified metrics for each pitch type, grouped by game_date and pitcher\n
avg_metrics = all_2018_stats_df.groupby(['game_date', 'pitcher', 'pitch_type']).agg({\n    'release_speed': 'mean',\n    'release_pos_x': 'mean',\n    'release_pos_z': 'mean',\n    'spin_dir': 'mean',\n    'vx0': 'mean',\n    'vy0': 'mean',\n    'vz0': 'mean',\n    'ax': 'mean',\n    'ay': 'mean',\n    'az': 'mean',\n    'effective_speed': 'mean',\n    'release_spin_rate': 'mean',\n    'release_extension': 'mean',\n    'release_pos_y': 'mean',\n    'spin_axis': 'mean'\n}).reset_index()\n\n# Merging total pitches and total pitches by type back\n
grouped_2018_df = total_pitches.merge(total_pitches_by_type, on=['game_date', 'pitcher'], how='left')\n
grouped_2018_df = final_df.merge(avg_metrics, on=['game_date', 'pitcher', 'pitch_type'], how='left')\n"

```

This information was saved to individual CSV's in order to save time and processing power.

```
In [45]: 1 """
          2 all_2018_stats_df.to_csv('all_2018_stats_df.csv')
          3 """
```

```
Out[45]: "\nall_2018_stats_df.to_csv('all_2018_stats_df.csv')\n"
```

This process was repeated for each season from 2008 - 2023

```
In [ ]: 1
```