

NBA Oyuncu Verilerinin Analizi

1st Mustafa Kemal Ekim
Computer Engineering Department
Yıldız Technical University
İstanbul, Turkey
kemal.ekim@std.yildiz.edu.tr

Abstract—Sporda veri madenciliği son yıllarda oldukça popüler hale gelmiştir. Veri madenciliği konuları, hesaplamalı ve istatistiksel olarak bu konuları derinlemesine araştırmamıza olanak sağlamaktadır. Çeşitli yöntemler kullanılarak takımların, oyuncuların hatta personel ekiplerinin veri analizi yapılarak son derece yararlı ve kullanışlı analizler yapılabilir. Bu makale, NBA oyuncularının 2021-2022 sezonundaki istatistiklerini kullanarak çıkarımlar ve belirli analizler yapılan projenin detaylarını sunmaktadır. Örneğin NBA liginde sezon sonunda en değerli oyuncu, en iyi savunma oyuncusu, en çok gelişme gösteren oyuncu gibi ödüller verilmektedir. Bu ödüller, jüriler tarafından belirlenen kriterler göz önüne alınarak puanlanmaktadır. Projeimizde ise veri setindeki değerler kullanılarak verinin işlenmesi sonucu ödülü gerçekten kimin hak ettiği bulunacaktır.

Index Terms—basketbol, spor, oyuncu, python

I. GİRİŞ

Spor analizi ve tahmini, dünyanın her yerinde popülerlik kazandı. Özellikle son yıllarda basketbol gibi daha iyi bilinen sporlar için ön plana çıktığını görüyoruz. Ulusal Basketbol Birliği veya NBA, Kuzey Amerika'nın en önde gelen erkek profesyonel basketbol ligidir. NBA'i oluşturan takımlar ve oyuncular için sayısız istatistik bulunmaktadır. Verilerin işlenmesi insanlar tarafından sağlarsa yapılan analizlerin çoğu duygulardan etkilenebilir. Toplanan NBA verileri üzerinde objektif olarak istatistiksel analizler yaparak, bu duygusal etkilerden bağımsız tahminler yapabilmek için veri madenciliğine başvururuz.

Özellikle duygudan bağımsız analiz ve karar verilmesi gereken durumlarda veri işleme konusu ön plana çıkmakta. Bu makale, NBA oyuncularının istatistiksel verileri üzerinde gerçekleştirilen çeşitli analizleri sunmakta ve bulgularını özetlemektedir. Belirli bir maçın sonucunu tahmin etmek veya en iyi savunmacıyı seçmek sınıflandırma problemleri arasına girer, oysa oynanan dakikalardan oyun başına skor gibi konuların tahmini bir lineer regresyon kullanılarak yanıtlanır.

II. İLGİLİ İŞLER

- Her spor dalı için takım oluştururken sporcu seçimi önemli bir işler. Bu seçimde sporcuların daha önceki karşılaşmalarda gösterdikleri performansın bilinmesi büyük önem taşımaktadır. Son yıllarda özellikle yurt dışında birçok spor ekibi yetenekli yeni oyuncular keşfetmek, mevcut oyuncuların eksikliklerini ortaya çıkarmak ve rakip durumunu analiz ettirmek amacıyla istatistikçi ve analist çalıştırmaktadır.

- Atletizm kategorilerindeki sporcuların fiziksel durumlarına ve performanslarına göre antrenman sezonunun etkinliğini hesaplamak için periyodik testler yapılır. Antrenörler, bilinen atletik verileri kullanarak sezon öncesi antrenmanlarını daha verimli hale getirebilirler veya test edilen aktivitelerin sporcular üzerindeki etkisini azaltılabilir. Ayrıca, belirli organizasyonlarda oyuncularının belirli koşullar altında nasıl performans göstereceklerini de tahmin edebilirler.

III. VERİLERİN ALINMASI

Bu proje için kullanılan veriler Basketball-Reference.com ve NBA.com sitelerinden sağlanmıştır. Belirtilen ilk site, her oyuncunun ve takımın tüm temel istatistiklerini içerir. Bu bilgiler oyun istatistiklerinden oyuncu maaşlarına kadar uzanır. [1] Verilere .CSV formatında kolayca erişebildik. Resmi olan NBA.com web sitesi ise bize her oyuncu için topa sahip olma süresi, şut mesafesi ve defans mesafesi gibi bir dizi daha gelişmiş istatistik sağladı. [2] Verileri aldığımız site basketbol oynamayan insanlara da hitap ettiği için istatistik sütunlarındaki yazılar açıklayıcı ve uzun bilgiler ile dolu idi.

```
1 data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 716 entries, 0 to 715
Data columns (total 28 columns):
#   Column
---  ---
0   FULL NAME
1   TEAM
2   POS
3   AGE
4   GP
5   MPG
6   MINMinutes PercentagePercentage of team minutes used by a player while he was on the floor
7   USGUsage RateUsage rate, a.k.a., usage percentage is an estimate of the percentage of team p
8   TOVTurnover RateA metric that estimates the number of turnovers a player commits per 100 poss
9   FTA
10  FT%
11  2PA
12  2P%
13  3PA
14  3P%
15  eFGEffective Shooting PercentageWith eFG%, three-point shots made are worth 50% more than tw
16  TSTrue Shooting PercentageTrue shooting percentage is a measure of shooting efficiency that
17  PPGPointsPoints per game.
18  RPGReboundsRebounds per game.
19  TRBTotal Rebound PercentageTotal rebound percentage is estimated percentage of available ret
20  APGAssistsAssists per game.
21  ASTAssist PercentageAssist percentage is an estimated percentage of teammate field goals a p
22  SPGStealsSteals per game.
23  BPGBlocksBlocks per game.
24  TOPGTurnoversTurnovers per game.
25  VIVersatility IndexVersatility index is a metric that measures a player's ability to produce
26  ORTGOffensive RatingIndividual offensive rating is the number of points produced by a player
27  DRGDefensive RatingIndividual defensive rating estimates how many points the player allowed
dtypes: float64(21), int64(4), object(3)
```

Pandas kütüphanesinin rename fonksiyonu sayesinde bu isimleri daha anlaşılır ve görüntü kirliliği yaratmayacak şekilde düzenledik.

```

1 data.rename(columns= {'eFG%': 'Effective Shooting Percentage', 'eFG%_True Shooting Percentage': 'True Shooting Percentage', 'TS%': 'True Shooting Percentage', 'TS%_True Shooting Percentage': 'True Shooting Percentage', 'TRB%': 'Total Rebound Percentage', 'TRB%_Total Rebound Percentage': 'Total Rebound Percentage', 'MIN%': 'Minutes Percentage', 'MIN%_Minutes Percentage': 'Minutes Percentage', 'USG%': 'Usage Rate', 'USG%_Usage Rate': 'Usage Rate', 'a.k.a.': 'a.k.a.', 'usage pe': 'Usage Rate', 'TO%': 'Turnover Rate', 'TO%_Turnover Rate': 'Turnover Rate', 'a metric that estimates th': 'a metric that estimates the turnover rate', 'APG': 'Assists per game', 'APG%': 'Assists per game', 'SPG': 'Steals per game', 'SPG%': 'Steals per game', 'BPG': 'Blocks per game', 'BPG%': 'Blocks per game', 'PPG': 'Points per game', 'PPG%': 'Points per game', 'RPG': 'Rebounds per game', 'RPG%': 'Rebounds per game', 'TOPG': 'Turnovers per game', 'TOPG%': 'Turnovers per game', 'V': 'Versatility Index', 'V%': 'Versatility Index', 'AST': 'Assist Percentage', 'AST%': 'Assist Percentage', 'ORTG': 'Offensive Rating', 'ORTG%': 'Offensive Rating', 'DRTG': 'Defensive Rating', 'DRTG%': 'Defensive Rating'})
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

```

IV. PYTHON

Python, genellikle karmaşık veri kümelerini düzene sokmak ve işlemek için kullanılan çeşitli avantajlara sahip, çok işlevli bir programlama dilidir. Python'un veri analizi için en iyi seçenek haline getiren bir dizi ayırt edici özelliği bulunmaktadır.

- Kolay Öğrenim: Python, basitlik ve okunabilirliğe odaklanan bir dil olmasına ek aynı zamanda veri analistleri için çok sayıda yararlı seçenek sunar. Python öğrenmeye yeni başlamış bir insan bile birkaç satır kod ile etkili çözümlere anında ulaşabilir.
- Esneklik: Python'un çok yönlülüğü, veri bilimcilerin onu popüler kılmasında etkili olmuştur. Veri modelleri oluşturma, veri setlerini işleme, makine öğrenimi algoritmaları geliştirme gibi görevleri kısa sürede tamamlamak için python kullanımı uygun olacaktır.
- Kütüphaneler: Python birçok kütüphaneye ücretsiz erişim sağlar. Bu kütüphanelerden bazıları Pandas, SciPy, scikit-learn, seaborn gibi kütüphanelerdir. Ayrıca bu kütüphanelerin sürekli geliştiğini ve ilgilendiğimiz konulara sağlam çözümler sunduğunu unutmamalıyız.
- Görselleştirme: İşlediğimiz verilerin rahat anlaşılabilmesi için görselleştirme çok önemlidir. Python, kullanıcılara basit kütüphaneler ile çok sayıda farklı görselleştirme seçeneği sunar. Bar ve dairesel grafikler, histogramlar, ısı haritaları gibi görselleştirme teknikleri veri analistleri tarafından sıkça kullanılan tekniklerdir.

Özetleyecek olursak, python yukarıda bahsedilen nedenlerle sektörde önemli bir yere sahiptir. Sürekli gelişmekte olması da uzun bir süre boyunca veri konularında adından söz ettireceğini bizlere göstermektedir.

V. TAHMIN

Projede bir model oluşturarak veriyi 80 train ve 20 test olacak şekilde ikiye böldük ve train kısmını eğitmeye başladık. Bu eğitimi 'LinearRegression' kütüphanesi ile gerçekleştirdik. Veri setinin eğitim ve test olarak bölünmesi önemli bir adımdır

çünkü doğru seçilmediği takdirde az öğrenme ya da çok öğrenme (ezberleme) gibi durumlar ortaya çıkabilir.

```

1 # 80 train, 20 test
2 from sklearn.model_selection import train_test_split
3 x_train, x_test, y_train, y_test = train_test_split(data[['PPG']], data[['APG']], test_size=0.2, random_state=2)

1 # Predicting APG from PPG
2 from sklearn.linear_model import LinearRegression
3 linear = LinearRegression()
4 linear.fit(x_train, y_train)
5 predicts = linear.predict(x_test)

```

Eğitim sonucunda iki istatistiğin birbirlerini tahmin etmesini 20'lik kısım üzerinden karşılaştırarak modelin doğruluk oranını gördük. İlişkisi yüksek olan sütunlar, örneğin oynanan dakika ve atılan sayı, daha yüksek bir doğruluk oranına sahipken korelasyonu düşük olan sütunlar, örneğin oyuncunun yaptığı sayı ve yaptığı asist, nispeten daha az bir doğruluk oranına sahip oldu. Aşağıdaki görsellerden bunu inceleyebilirsiniz.

```

1 # success rate
2 linear_score = linear.score(x_test, y_test)
3 print("Score: ", linear_score)
4
5 # Since points and assists are different concepts, the accuracy rate is low.

```

Score: 0.4947593819968187

```

1 x_train2, x_test2, y_train2, y_test2 = train_test_split(data[['MPG']], data[['PPG']]
2 linear2 = LinearRegression()
3 linear2.fit(x_train2, y_train2)
4 predicts = linear2.predict(x_test2)
5 print("Score: ", linear2.score(x_test2, y_test2))
6
7 # The more time a player has, the more chances to score.

```

Score: 0.738672950115288

VI. KÜMELEME

Kümeleme, birbirine benzeyen nesnelerin aynı kümede yer almasını sağlayarak veri setini istenen sayıda kümeye ayırma problemidir. Bu projede 'K-means clustering' adı verilen gözetimsiz öğrenme algoritmasını kullandık. Küme sayısını belirlemek algoritma açısından oldukça önemlidir. Eğer k sayısı az sayıda ise birbirine benzemeyen nesneler aynı kümeye düşecektir, eğer k sayısı çok fazla ise kümeleme işlemi anlamsız olacaktır çünkü neredeyse her nesnenin kendine ait kümesi olacaktır, bu istenmeyen bir durumdur. K sayısını seçmek için 'The Elbow Method' ve 'The Silhouette Method' gibi yöntemler kullanılabilir, bu makalede detayları işlenmeyecektir. Projede sklearn.cluster kütüphanesinin KMeans fonksiyonu kullanılarak oyuncular 5 farklı kümeye dağıtılmıştır. Sonrasında PCA (Principal Component Analysis) kullanılarak her oyuncu bir kümeye atanmıştır.

```

1 #Cluster of players (k-means)
2 from sklearn.cluster import KMeans
3 model = KMeans(n_clusters=5, random_state=1) # KMeans model (5 clusters)
4 efficient_cols = data.get_numeric_data().dropna(axis=1) # get pure data
5 model.fit(efficient_cols) # training
6 labels = model.labels_
7 labels

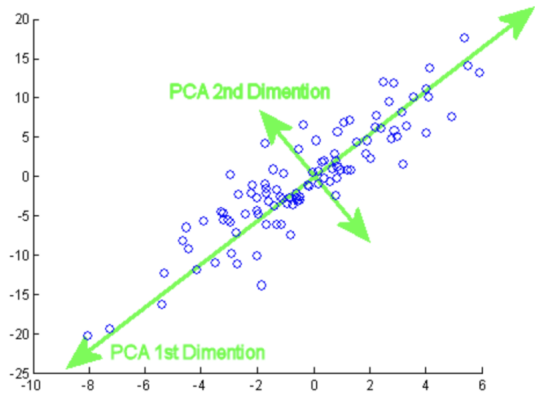
```

```

array([2, 2, 1, 4, 2, 3, 4, 3, 2, 0, 4, 4, 2, 1, 4, 3, 3, 3, 4, 4, 4, 4,
2, 4, 2, 4, 0, 0, 1, 0, 3, 3, 0, 4, 3, 1, 1, 3, 4, 4, 0, 0, 4, 2,
2, 3, 4, 0, 4, 4, 4, 0, 3, 0, 0, 0, 0, 2, 4, 3, 3, 4, 4, 1, 0,
2, 4, 0, 4, 0, 1, 2, 2, 2, 0, 4, 2, 4, 4, 4, 4, 0, 0, 2, 1, 4, 4,
4, 1, 4, 4, 3, 0, 3, 4, 1, 4, 3, 0, 4, 2, 4, 4, 2, 4, 4, 0, 4,
4, 4, 4, 4, 4, 2, 2, 4, 3, 0, 0, 2, 0, 4, 3, 3, 4, 4, 4, 4, 0,
0, 4, 0, 4, 3, 4, 4, 1, 0, 0, 3, 2, 4, 0, 4, 0, 4, 4, 1, 4, 4, 2,
4, 2, 0, 4, 0, 1, 3, 2, 4, 4, 4, 4, 4, 4, 4, 4, 0, 0, 2, 4, 4,
1, 4, 1, 4, 0, 4, 0, 0, 1, 4, 4, 4, 0, 0, 4, 4, 4, 3, 4, 4, 0,
0, 0, 4, 0, 3, 1, 4, 4, 4, 4, 0, 4, 4, 2, 3, 4, 4, 1, 4, 4, 4,
0, 2, 0, 2, 1, 4, 4, 2, 4, 4, 2, 3, 3, 2, 4, 0, 0, 0, 3, 0, 2, 0,
0, 0, 4, 0, 2, 0, 4, 2, 3, 3, 0, 0, 4, 2, 0, 3, 4, 4, 1, 2, 0, 0,
4, 4, 2, 2, 4, 4, 4, 4, 2, 1, 3, 0, 4, 0, 4, 4, 4, 4, 4, 0, 4,
1, 0, 0, 0, 4, 4, 4, 3, 2, 4, 4, 4, 0, 4, 3, 4, 4, 4, 2, 4, 4, 3,
0, 4, 4, 0, 1, 2, 4, 1, 0, 0, 0, 4, 4, 4, 1, 1, 4, 0, 4, 4, 4,
4, 4, 3, 4, 0, 4, 4, 1, 4, 0, 4, 4, 1, 0, 4, 4, 0, 2, 4, 4, 4,
0, 2, 4, 4, 0, 4, 4, 3, 4, 0, 4, 4, 4, 3, 0, 4, 4, 4, 4, 0, 4,
4, 4, 4, 4, 2, 3, 4, 0, 4, 0, 0, 4, 1, 4, 0, 0, 0, 2, 0, 4, 4, 3,
0, 4, 2, 4, 0, 4, 4, 3, 3, 4, 0, 0, 0, 0, 4, 2, 3, 4, 3, 0, 2, 0,
0, 4, 4, 3, 4, 4, 0, 1, 4, 4, 4, 3, 2, 0, 0, 2, 0, 2, 4, 0, 4,
4, 0, 3, 3, 4, 0, 1, 4, 4, 4, 3, 0, 3, 1, 1, 3, 4, 4, 4, 4, 0, 4,
4, 1, 4, 4, 2, 4, 4, 3, 4, 4, 4, 0, 1, 0, 0, 4, 0, 4, 4, 0, 2,
4, 3, 4, 0, 4, 0, 2, 0, 0, 4, 4, 3, 4, 0, 2, 4, 4, 4, 0, 4, 0,
4, 3, 4, 3, 4, 4, 4, 4, 0, 2, 2, 0, 4, 4, 4, 4, 4, 2, 2, 0, 4,
3, 3, 4, 0, 3, 2, 0, 4, 2, 4, 2, 4, 0, 0, 0, 4, 4, 3, 4, 1, 0, 0,
4, 4, 4, 2, 4, 0, 0, 0, 3, 0, 4, 4, 4, 2, 0, 4, 4, 0, 3, 3, 0, 3,
4, 2, 0, 4, 4, 4, 4, 4, 4, 2, 4, 4, 4, 4, 2, 4, 3, 0, 1, 4, 4, 4,
3, 4, 4, 4, 4, 4, 3, 0, 0, 0, 4, 0, 4, 4, 4, 4, 4, 4, 4, 4, 0,
0, 2, 3, 2, 4, 0, 2, 1, 0, 0, 4, 4, 4, 0, 4, 4, 2, 4, 4, 4, 3,
0, 4, 4, 4, 0, 4, 0, 4, 2, 0, 1, 3, 0, 4, 4, 0, 1, 4, 4, 2, 3, 3,
3, 1, 0, 1, 0, 4, 3, 0, 4, 4, 4, 4, 4, 0, 3, 4, 4, 4, 4, 0, 4, 4,
1, 3, 2, 0, 2, 4, 0, 1, 4, 2, 0, 0, 0, 0, 4, 4, 4, 4, 4, 4, 4,
1, 4, 0, 4, 4, 4, 4, 4, 1, 0, 4, 2], dtype=int32)

```

PCA'dan kısaca bahsedecek olursak, principal component analysis yani temel bileşen analizi, çok boyutlu bir uzaydaki verilerin, varyansı en üst düzeye çıkaracak şekilde daha düşük boyutlu bir alana izdüşümü bulma yöntemidir. Uzaydaki bir dizi nokta için, tüm noktaların en düşük ortalama mesafesine sahip "en uygun çizgi" seçilir. [3]



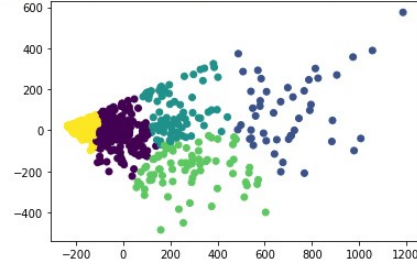
[4] Bu işlemin görselleştirilmesi ve bazı örnekleri aşağıdaki fotoğraflardan inceleyebilirsiniz.

```

1 # plot clusters
2 from sklearn.decomposition import PCA
3 pca_2 = PCA(2)
4 plot_cols = pca_2.fit_transform(efficient_cols)
5 plt.scatter(x=plot_cols[:,0], y=plot_cols[:,1], c=labels)

```

<matplotlib.collections.PathCollection at 0x7f3f32d527d0>



```

1 LeBron = efficient_cols.loc[ data['FULL NAME'] == 'LeBron James', :]
2
3 Jokic = efficient_cols.loc[ data['FULL NAME'] == 'Nikola Jokic', :]
4
5 Mikal = efficient_cols.loc[ data['FULL NAME'] == 'Mikal Bridges', :]

```

```

1 LeBron_list = LeBron.values.tolist()
2 Jokic_list = Jokic.values.tolist()
3 Mikal_list = Mikal.values.tolist()
4
5 print("LeBron's Cluster: ", model.predict(LeBron_list))
6 print("Jokic's Cluster: ", model.predict(Jokic_list))
7 print("Mikal's Cluster: ", model.predict(Mikal_list))

```

```

LeBron's Cluster: [1]
Jokic's Cluster: [1]
Mikal's Cluster: [2]

```

VII. KORELASYON

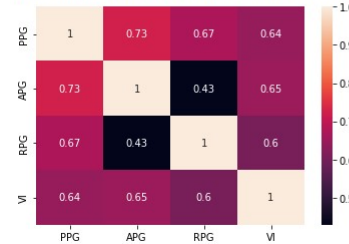
Korelasyon analizi, iki değişken arasındaki ilişkiyi hesaplamak için kullanılan istatistiksel bir yöntemdir. Bir sütundaki verilerin, diğer sütundaki verilerin değişimiyle ne kadar ilişkide olduğunu korelasyon katsayısı ile inceleriz. [5] Aşağıdaki örnekte sayı, asist, ribaund ve versatility index'in korelasyonunu ve ısı haritasını görebilirsiniz.

```

1 import seaborn as sns
2 #Heat map
3 # VI = measures a player's ability to produce in points, assists, and rebounds
4 correlation = data[['PPG', 'APG', 'RPG', 'VI']].corr()
5 sns.heatmap(correlation, annot=True)

```

<matplotlib.axes._subplots.AxesSubplot at 0x7f3f39bfd990>



VIII. MVP ÖDÜLÜ

Sezonun en değerli oyuncusu ödülü jürilerin verdiği oylar sonucunda belirlendiği için duygu durumlarından etkilenmesi oldukça muhtemeldir. Veri madenciliği sayesinde önemli kriterler işlenerek objektif bir biçimde sezonun en değerli oyuncusu belirlenebilir. Projede kriterleri kendim seçtim ve oranladım. Bu kriterler oynanan maç +66 olmalı, en az 25 sayı

ortalaması, en az 4 asist ortalaması, en az 4 ribaund ortalaması ve şut yüzdesinin 40'tan fazla olması oldu. Bunun sonucunda 716 kişiden sadece 7 kişi MVP yarışına kalabildi.

```
3 mvp_candidates[['FULL_NAME', 'GP', 'PPG', 'APG', 'RPG', '2P%']]
```

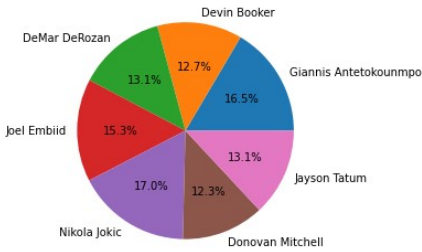
	FULL_NAME	GP	PPG	APG	RPG	2P%
13	Giannis Antetokounmpo	67	29.9	5.8	11.6	0.616
64	Devin Booker	68	26.8	4.8	5.0	0.508
150	DeMar DeRozan	76	27.9	4.9	5.2	0.520
184	Joel Embiid	68	30.6	4.2	11.7	0.529
343	Nikola Jokic	74	27.1	7.9	13.8	0.652
453	Donovan Mitchell	67	25.9	5.3	4.2	0.533
623	Jayson Tatum	76	26.9	4.4	8.0	0.524

Kalan 7 adayın kriterlerdeki ortalaması alındı ve belirlenen yüzdelere göre her adayın bir puanı ortaya çıktı.

	FULL_NAME	GP	PPG	APG	RPG	eFG%	MVP_Points
343	Nikola Jokic	74	27.1	7.9	13.8	0.620	12.974156
13	Giannis Antetokounmpo	67	29.9	5.8	11.6	0.582	12.573790
184	Joel Embiid	68	30.6	4.2	11.7	0.534	11.640985
150	DeMar DeRozan	76	27.9	4.9	5.2	0.521	9.988546
623	Jayson Tatum	76	26.9	4.4	8.0	0.526	9.943546
64	Devin Booker	68	26.8	4.8	5.0	0.530	9.650985
453	Donovan Mitchell	67	25.9	5.3	4.2	0.533	9.378790

Hesaplanan puanlara göre Nikola Jokic MVP ödülünün sahibi oldu. Gerçekte de ödülü kazanan o olduğundan jüriye duyugudan bağımsız puanlar verdiğini söyleyebiliriz. Aşağıdaki dairesel grafikte hangi oyuncunun ödülü yüzde kaç hak ettiğini görebilirsiniz.

```
1 # Ağırlıklı kriterlere göre mvp tahminleri ve pie chart
2 mvp_points
3
4
5
6 fig1, ax1 = plt.subplots()
7 ax1.pie(mvp_candidates['MVP_Points'], labels=mvp_candidates['FULL_NAME'],
8 ax1.axis('equal')
9 plt.show()
```



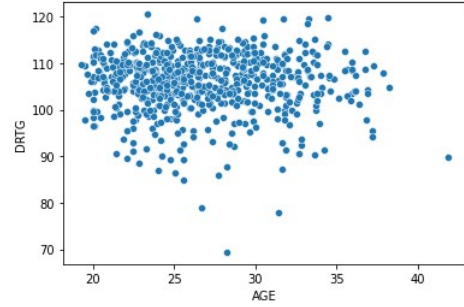
IX. AYKIRI VERİLERİN AYRIŞTIRILMASI

Özetle, aykırı veri, veri kümesindeki diğer gözlemlerden önemli ölçüde farklı olan herhangi bir veri alanıdır. Aykırı verilerin bulunduğu kümedeki veya toplama verilerindeki

diğer verilerden farklı davranır ve hatalara neden olur.

Aykırı verileri temizlerken; Kutu grafiği kullanma, 5 sayı özeti ve standart sapma ile uç değerleri belirleme gibi yöntemler kullanılır. Aykırı değerlerin çıkarılması, ortalama ile doldurulması veya silme işlemi ile biter. [6]

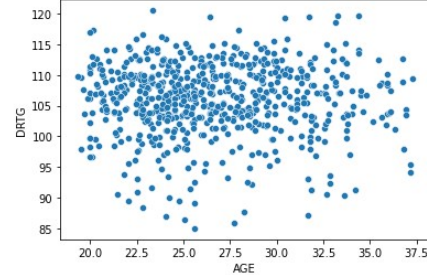
```
1 # MPG and PPG correlation (line graph)
2 ax = sns.scatterplot(x="AGE", y="DRTG", data=data)
3 ax.set_title("Age and Defence Correlation")
```



Yukarıdaki fotoğrafta yaş ve defansif güç arasındaki bağlantıyı noktalar şeklinde görüyoruz. Aşağıdaki işlemde uç noktaları veri setimizden çıkarıp işlemlere devam ediyoruz.

```
1 wo_outliers = data[(data['AGE'] > 37.5) | (data['DRTG'] < 85)]
2 wo_outliers = pd.concat([data, wo_outliers]).drop_duplicates(keep = False)
3
4 sns.scatterplot(x="AGE", y="DRTG", data=wo_outliers)
```

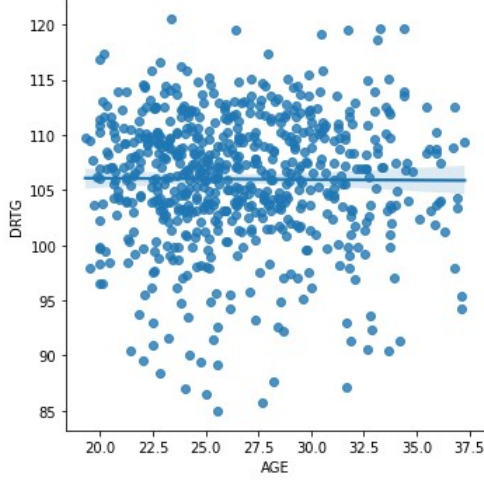
<matplotlib.axes._subplots.AxesSubplot at 0x7f3f2feb5650>



Defansif güç özelliği anormal derecede düşük olan oyuncular ve yaşı çok büyük oyuncular veri setimizden çıkarıldığında ortada daha tutarlı bir veri kaldığını görüyoruz.


```
1 sns.lmplot(x="AGE", y="DRTG", data=wo_outliers)
```

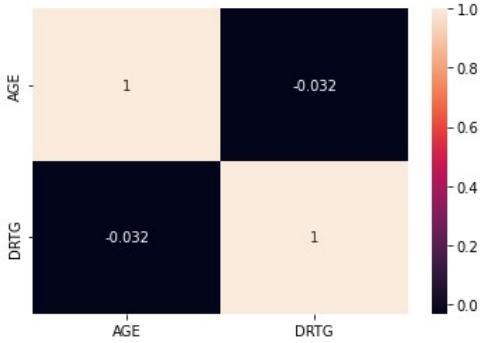
<seaborn.axisgrid.FacetGrid at 0x7f3f2ff68dd0>



Bir dağılım grafiğine en uygun çizgiyi eklemek bazı durumlarda yararlı olabilir ancak elimizde çok fazla sayıda oyuncu verisi olduğu için bu örnekte efektif olduğunu söyleyemeyiz.

```
1 correlation_age_drtg = data[['AGE', 'DRTG']].corr()  
2 sns.heatmap(correlation_age_drtg, annot=True)
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f3f30ba5b90>



İncelediğimiz örnekte korelasyonun eksi 0.032 seviyesinde çıktığını gördük. Bu gayet anlaşılabilir bir durum çünkü bir oyuncunun yaşı ilerledikçe defans konusunda güçlük çekmesi normaldir çünkü genç oyunculara karşı yavaş kalabilir, daha çabuk yorulabilir. Eğer örneğimizde yaş - hücum ikilisinin ilişkisi inceleysek muhtemelen artı seviyede bir korelasyon sayısı karşımıza çıkacaktı çünkü oyuncu ne kadar fazla süre basketbol oynar veya izlerse oyunu o kadar iyi tanır ve hücumda ne yapması gerektiğini daha iyi bilir.

X. SONUÇ

Günümüz dünyasında her alanda bir çok veri üretilmektedir ve artık insanların bu verileri kâğıt kalem ile analiz etmesi mümkün olmamaya başlamıştır. Eğitim, ticaret, bilim, finans, sağlık ve spor gibi bir çok alanda işe yarar verileri depolama ve işleme günümüzde revaçtadır ve uzun bir süre boyunca popüler kalmaya devam edecek gibi durmaktadır.

Hazırlanan proje ve bu yazıda, NBA istatistiksel verileri

üzerinde gerçekleştirilen birkaç analiz sunduk ve sonuçları özetledik. Bu makale basketbol taraftarlarına bilgileri daha iyi tanıtacak, bu arada basketbolla ilgilenmeyen insanlar için basit bir tanıtım olmaktadır. Python dili sayesinde veri setimiz üzerinde tahmin ve doğruluk oranı, kümeleme, korelasyon, aykırı veriler gibi konular projede işlenmiştir.

REFERENCES

- [1] <https://www.basketball-reference.com/>
- [2] <https://www.nba.com/>
- [3] <https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.PCA.html>
- [4] <https://programmatically.com/principal-components-analysis-explained-for-dummies/>
- [5] <https://tr.wikipedia.org/wiki/Korelasyon>
- [6] <https://machinelearningmastery.com/how-to-use-statistics-to-identify-outliers-in-data/>