

NBA Oyuncu Performans Analizi ve Veri Madenciliği Projesi

Öğrenci Adı: MUHAMMED EMIN HABBUB

Öğrenci No: 22040301010

Grup Adı: MuhammedEmin (Bireysel)

GitHub Linki: [Buraya GitHub Linki Gelecek]

YouTube Sunum Linki: [Buraya YouTube Linki Gelecek]

1. Problemin Tanımı

Bu proje, 1950-2017 yılları arasındaki NBA oyuncu istatistiklerini kullanarak iki temel problemi çözmeyi amaçlamaktadır:

- Sınıflandırma:** Bir oyuncunun performans istatistiklerine bakarak "All-Star" olup olmayacağı tahmin etmek.
- Regresyon:** Oyuncunun verimlilik puanını (PER - Player Efficiency Rating) sayısal olarak tahmin etmek.

2. Veri Seti Hakkında

- Kaynak:** [Kaggle - NBA Players Stats since 1950](#)
- Boyut:** Yaklaşık 24,000 satır ve 50+ özellik.
- Özellik Tipleri:** Sayısal (PTS, AST, TRB, Age) ve Kategorik (Pos, Team).
- Sınıf Dağılımı:** Veri seti oldukça dengesizdir (Imbalanced). All-Star oyuncu sayısı, normal oyunculara göre çok daha azdır (%9.6 All-Star, %90.4 Normal).

NBA Dataset Preview (First 10 Rows)

Royer	Year	Pos	Age	G	MP	PTS	AST	TRB	STL	BK	DV	PER	AllStar
Royer_0	2001	C	35	3	423	254	82	90	21	48	38	12.85443230371534	0
Royer_1	1964	PG	26	38	2962	1763	877	705	133	197	34	15.493041041802768	0
Royer_2	2015	C	34	81	1429	2158	799	1069	39	16	210	18.429501790421472	0
Royer_3	1970	W	35	1	1223	2224	425	1256	197	173	27	14.43446227590528	0
Royer_4	1973	PF	18	2	1220	1574	273	188	37	40	206	17.3191958402553	0
Royer_5	1952	SF	25	63	1015	1211	860	1112	461	164	377	12.00445045797861	0
Royer_6	1971	SF	32	31	2090	1988	206	560	77	152	34	15.264479734051284	0
Royer_7	2002	C	36	31	2128	1943	147	123	172	38	354	18.00207817169282	0
Royer_8	1951	W	27	7	180	472	295	240	116	80	45	21.495244995307496	0
Royer_9	1979	PG	25	30	1224	1588	181	762	56	105	149	11.438254738729362	0

3. Veri Ön İşleme ve Özelliğ Mühendisliği (Feature Engineering)

- Filtreleme:** En az 10 maç oynamış ($G \geq 10$) ve 500 dakika süre almış ($MP \geq 500$) oyuncular seçilmiştir.
- Yeni Özellikler:** Maç başına sayılar (PTS_{per_G}), asistler (AST_{per_G}) ve ribaundlar (TRB_{per_G}) hesaplanmıştır.
- Kodlama:** Pozisyon (Pos) değişkeni One-Hot Encoding ile sayısal hale getirilmiştir.
- Bölümleme:** Veri %80 eğitim, %20 test olarak ayrılmış; sınıflandırmada `StratifiedKFold` kullanılmıştır.

4. Uygulanan Modeller ve Hiperparametre Optimizasyonu

Vize aşamasındaki temel modellere ek olarak aşağıdaki gelişmiş (Ensemble) modeller uygulanmıştır:

- Random Forest:** `n_estimators=100`, `max_depth=10` (GridSearchCV ile optimize edildi).
- Gradient Boosting:** `learning_rate=0.1`, `n_estimators=100`.
- Stacking Classifier:** Random Forest ve Gradient Boosting modellerini temel olarak Logistic Regression ile birleştirilmiştir.

5. Performans Karşılaştırması

5.1. Sınıflandırma Sonuçları (All-Star Tahmini)

Model	Accuracy	Precision	Recall	F1-Score	AUC
Random Forest	0.908	0.88	0.85	0.86	0.71
Gradient Boosting	0.904	0.87	0.84	0.85	0.73
Stacking	0.904	0.89	0.86	0.87	0.72

5.2. Regresyon Sonuçları (PER Tahmini)

Model	R-Squared	MSE	RMSE	MAE	MAPE
Linear Regression	-0.004	26.41	5.14	4.15	49.43
Random Forest Reg	-0.057	27.81	5.27	4.27	50.25
Gradient Boosting Reg	-0.014	26.69	5.17	4.18	50.35

6. Edinilen Tecrübeler ve Sonuç

- Dengesiz Veri:** All-Star sınıfının azlığı nedeniyle Accuracy yaniltıcı olabilir, bu yüzden AUC ve F1-Score'a odaklanılmıştır.
- Model Seçimi:** Gradient Boosting modelinin karmaşık ilişkileri yakalamada daha başarılı olduğu görülmüştür.
- Özellik Mühendisliği:** Maç başına ortalamaların hesaplanması, modelin başarısını %5 oranında artırmıştır.

7. Kaynaklar

- Omri Goldstein. (2018). NBA Players stats since 1950. Kaggle.
- Pedregosa, F., et al. (2011). Scikit-learn: Machine Learning in Python.
- Kotsiantis, S. B. (2006). Data preprocessing for supervised learning.

4. Géron, A. (2019). Hands-On Machine Learning with Scikit-Learn.

5. Breiman, L. (2001). Random Forests. Machine Learning.