

# Maç Sonucu Tahmini İçin Veri Seti Oluşturma

Ensar Akbaş  
Bilişim Sistemleri Mühendisliği  
Kocaeli Üniversitesi  
ensar.akbas@gmail.com

**Özet**—Bu projede, futbol maç sonucu tahmin modelleri için kapsamlı bir veri seti oluşturulmuştur. Veri, Maçkolik sitesinden web scraping yöntemiyle toplanmıştır. Liverpool, Manchester City, Manchester United, Arsenal, Chelsea ve Tottenham takımlarına ait bilgiler derlenmiştir. Ham verilerde yer alan temel istatistiklere ek olarak, özellik mühendisliği (feature engineering) teknikleriyle yeni değişkenler üretilmiştir. Eksik veriler uygun yöntemlerle doldurulmuş, korelasyon analizi ve görselleştirmeler ile veri yapısı detaylı olarak incelenmiştir. Proje, maç istatistiklerinden öğrenebilen modeller için kaliteli, temiz ve anlamlı bir veri seti oluşturmayı hedeflemektedir.

**Keywords**—Web scraping, Feature engineering, Veri seti, Veri ön işleme

## I. GİRİŞ

Futbol, yalnızca bir spor dalı olmanın ötesinde, büyük bir veri kaynağı ve analiz alanı haline gelmiştir. Bu projede amaç, geçmiş futbol maçlarından elde edilen istatistiksel verileri kullanarak ileriye dönük tahminler yapılabilecek nitelikte bir veri seti oluşturmaktır. Özellikle İngiltere Premier Lig'in önde gelen takımlarına ait uzun dönemli maç verileri incelenmiş, her bir maçın istatistiksel özellikleri detaylı şekilde toplanarak analiz edilebilir bir yapıya dönüştürülmüştür. Proje kapsamında web scraping, veri temizleme, veri ön işleme ve feature engineering işlemleri gerçekleştirilmiştir.

## II. VERİ TOPLAMA

Veriler, Maçkolik web sitesinden Python programlama dili kullanılarak çekilmiştir. Kodlama işlemleri PyCharm editörü üzerinde yürütülmüş, verileri çekmek için Selenium kütüphanesi tercih edilmiştir. Ayrıca verilerin işlenmesi ve analizinde pandas, numpy ve görselleştirme için matplotlib kütüphanelerinden yararlanılmıştır. Süreç boyunca altı farklı takımın (Liverpool, Manchester City, Manchester United, Chelsea, Arsenal, Tottenham) her biri için yaklaşık 700 maçlık veri toplanmıştır. Bu işlemler sırasında Google Chrome Driver kullanılmış ve veri çekimi sırasında Maçkolik platformu herhangi bir bot koruması veya erişim engeli uygulamadığı için süreç kesintisiz ve verimli bir şekilde tamamlanmıştır.

## III. VERİ SETİ

Projenin veri seti, 2013-2014 sezonu ile 2024-2025 sezonu arasındaki Premier Lig maçlarını kapsamaktadır. Liverpool, Manchester City, Manchester United, Chelsea, Arsenal ve Tottenham olmak üzere 6 farklı takım için oluşturulan bu veri setinde, her takım adına yaklaşık 700 maç verisi bulunmaktadır. Bu kapsamda, toplam 32 özellik (feature) yer almaktadır. Özelliklerin bir bölümü doğrudan maç istatistiklerinden elde edilirken; bazıları ise özellik mühendisliği (feature engineering) süreçleriyle türetilmiştir. Bu yapı, takımların geçmiş performanslarını daha anlamlı şekilde modelleyebilmek için tasarlanmıştır.

## A. Features (Özellikler)

Tarih	Rakip Takım	Takım ID
Is_Home	Sonuç	Gol
Rakip Gol	Topla Oynama(%)	Şut
İsabetli Şut	Başarılı Pas	Pas Başarısı(%)
Korner	Faul	Ofsayt
Rakip Topla Oynama (%)	Rakip Şut	Rakip İsabetli Şut
Rakip Başarılı Pas	Rakip Pas Başarısı(%)	Rakip Korner
Rakip Faul	Rakip Ofsayt	

Tablo 1. Maçkolik üzerinden çekilen özellikler

- Takım ID:** Takıma özgü kimlik numarasıdır. (Örneğin: Liverpool=1, Manchester City=2)
- Is\_Home:** Maçın ev sahibi olarak oynanıp oynanmadığını belirtir. (Ev Sahini=1, Deplasman=0)
- Sonuç:** Galibiyet: 1, Beraberlik: 0, Mağlubiyet: -1 olarak tutulmaktadır.

## B. Feature Engineering (Özellik Mühendisliği)

Şut Verimliliği (İsabetli Şut/Şut)	Gol Farkı	Sezon
Ay	Haftanın Günü	Son 5 Maçın Gol Ortalaması
Son 5 Maçın Galibiyet Oranı		

Tablo 2. Feature engineering sonucu oluşan yeni özellikler

- Sezon/Ay:** Maçın hangi yıl ve ayda yapıldığını ifade eder. (Örneğin: Sezon=2024, Ay=4)
- Haftanın Günü:** Maçın hangi gün yapıldığını ifade eder. (Örneğin: 1=Pazartesi, 2=Salı)

Bu veri seti, sonraki aşamalarda analiz, görselleştirme ve makine öğrenmesi modelleri için temel kaynak olarak kullanılacaktır.

## IV. VERİ ÖN İŞLEME VE TEMİZLEME

Toplanan ham veriler, analiz ve modelleme süreçlerinde kullanılmadan önce çeşitli temizleme ve ön işleme adımlarından geçirilmiştir. Bazı maçlarda istatistiklerin eksik olması nedeniyle, ilgili sütunlarda boş (null) değerler tespit edilmiştir. Bu durum özellikle "Topla Oynama", "Pas Başarısı", "Başarılı Pas" gibi istatistiklerde görülmüştür. Bu eksik verileri doldurmak amacıyla iki farklı yaklaşım kullanılmıştır:

**Sütun Ortalamasıyla Doldurma:** Eksik değerler, ilgili sütunun genel ortalaması ile doldurularak, verinin genel yapısı korunmuştur.

**Önceki Verilere Yakınlaştırma:** Bazı durumlarda, eksik değerler bir önceki veya benzer maçlardaki istatistikler göz önünde bulundurularak manuel olarak doldurulmuştur. Böylece ani sapmaların önüne geçilmiştir.

Ayrıca, verilerin bütünlüğünü korumak adına tarih biçimleri, takım ID'leri ve sonuç değerleri gibi temel alanlarda format düzeltmeleri ve tutarlılık kontrolleri yapılmıştır. Bu adımlar, modelin güvenilir sonuçlar üretebilmesi için temel bir hazırlık sürecidir.

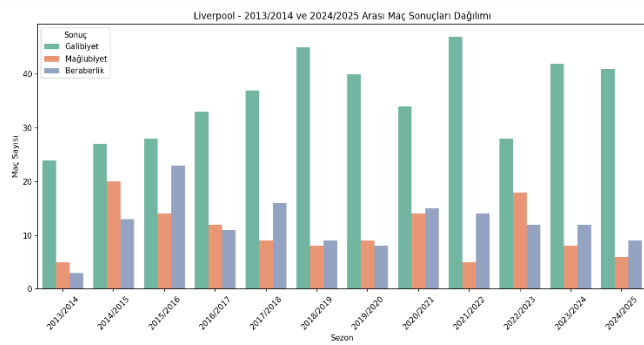
## V. VERİ GÖRSELLEŞTİRME

Veri analizi sürecinde, istatistiksel içgörülerin daha anlaşılır ve yorumlanabilir hale gelmesi için veri görselleştirme önemli bir adımdır. Bu projede, görselleştirme işlemleri hem genel eğilimleri ortaya koymak hem de değişkenler arasındaki ilişkileri incelemek amacıyla gerçekleştirilmiştir. Özellikle Liverpool takımı üzerinde yapılan analizlerle örneklemeler sunulmuştur.

### Liverpool Üzerinden Görselleştirme Örnekleri

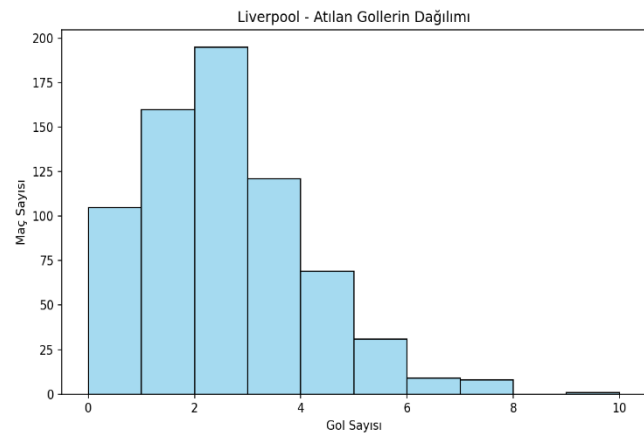
#### A. Sezonlara Göre Maç Sonuçları Dağılımı

Bu grafik ile her sezon boyunca kazanılan, kaybedilen ve berabere biten maçların dağılımı görselleştirilmiştir. Böylece takımın dönemsel performans dalgalanmaları net bir şekilde gözlemlenmiştir.



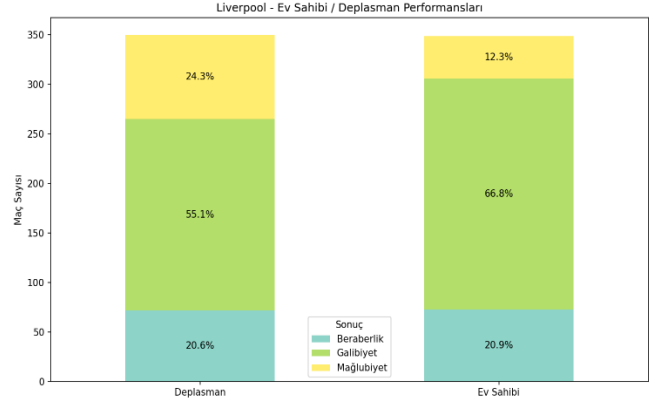
#### B. Atılan Gollerin Dağılımı

Takımın maç başına attığı gollerin frekansı incelenmiş, hangi gol aralığında en çok maç oynandığı grafiklerle gösterilmiştir.



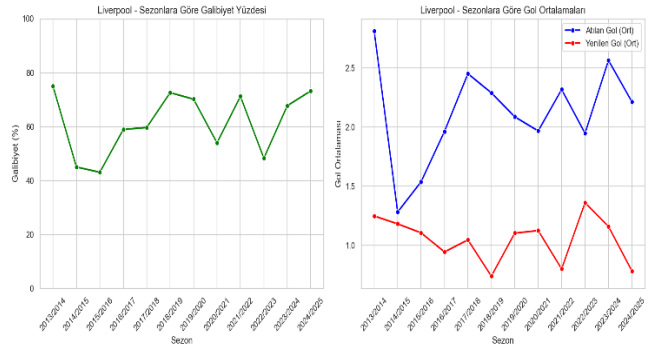
#### C. Ev Sahibi / Deplasman Performansları

Takımın iç saha ve dış saha maçlarındaki galibiyet oranları karşılaştırılarak performans farkları ortaya konmuştur.



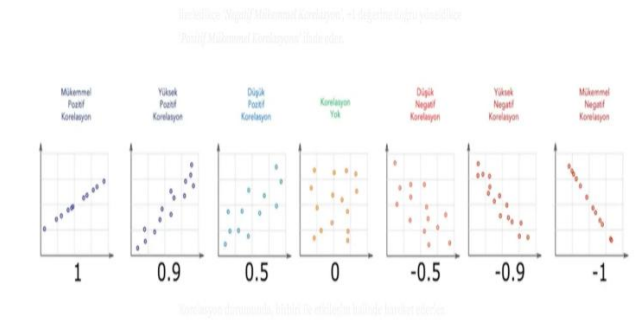
#### D. Sezonlara Göre Galibiyet Yüzdesi ve Gol Ortalamaları

Her sezonun galibiyet yüzdesi ile maç başına düşen gol ortalaması birlikte incelenerek takımın genel başarımı değerlendirilmiştir.



### Korelasyon Analizi

Korelasyon analizi, iki değişken arasındaki doğrusal ilişkinin yönünü ve gücünü ölçen istatistiksel bir yöntemdir. Bu projede, maç istatistikleri arasındaki ilişkileri incelemek amacıyla korelasyon analizinden yararlanılmış ve sonuçlar aşağıdaki başlıklar altında görselleştirilmiştir.



#### E. Korelasyon Matrisi

Tüm sayısal değişkenler arasındaki korelasyon katsayıları ısı haritası ile gösterilerek, hangi değişkenlerin birbirini daha fazla etkilediği görsel olarak sunulmuştur.

