

EDA ve Preprocessing Raporu

1. Veri Seti Bilgisi

- Satır/sütun sayısı: 2235 x 13
- Sayısal değişkenler: `Yas`, `TedaviSuresi`, `UygulamaSuresi`
- Kategorik değişkenler: `Cinsiyet`, `KanGrubu`, `Uyruk`, `Bolum`, `Alerji`, `Tanilar`, `TedaviAdi`, `UygulamaYerleri`

Gözlem: Veri setinde çoğunlukla Türk hastalara ait bilgiler bulunmakta ve sayısal değişkenler tedavi ile ilgili süreleri temsil ediyor.

2. Eksik Değer Analizi

- `Cinsiyet`: %7.5 eksik
- `KanGrubu`: %30 eksik
- `KronikHastalik`: %27 eksik
- `Alerji`: %42 eksik
- `Tanilar`: %3.3 eksik
- `UygulamaYerleri`: %9.8 eksik

Gözlem: En çok eksik değer `Alerji` değişkeninde bulunuyor, veri ön işleme sırasında uygun stratejilerle dolduruldu.

3. Aykırı Değerler

- `Yas`: 41 aykırı değer → Winsorize ile sınırlandırıldı
- `TedaviSuresi`: 565 aykırı değer → clip yöntemiyle sınırlandırıldı
- `UygulamaSuresi`: 12 aykırı değer → clip yöntemiyle sınırlandırıldı

Gözlem: Aykırı değerler, veri dağılımını bozmayacak şekilde düzeltilmiştir.

4. Ön İşleme Adımları

- `TedaviSuresi` ve `UygulamaSuresi` sayısala çevrildi
- Eksik değerler uygun stratejilerle dolduruldu (mod, 'Yok' vb.)
- Metin verilerinden çıkarılan sayısal özellikler:
 - `KronikHastalik_sayisi`
 - `Alerji_sayisi`
 - `Tanilar_sayisi`
- NLP Feature Engineering: Tanılar ve kronik hastalık kelime sıklığı çıkarıldı
- Kategorik değişkenler OneHotEncode edildi
- Sayısal değişkenler StandardScaler ile standartlaştırıldı

Gözlem: Bu adımlar ile veri modellemeye hazır hâle getirildi.

5. EDA ve Görselleştirmeler

- Sayısal deęiřkenlerin histogram ve boxplotları
- Korelasyon heatmap: `TedaviSuresi` ile en iliřkili deęiřkenler grselleřtirildi
- Encode edilmiř kategorik deęiřkenlerin frekans grafikleri
- NLP zetleri: En sık geen 10 tanı ve kronik hastalık kelimesi
- Tedavi Sresi iliřkili grselleřtirmeler:
 - Cinsiyet vs `TedaviSuresi`
 - Yař vs `TedaviSuresi`

Gzlem: Tedavi sresi ile bazı kronik hastalık ve yař deęiřkenleri arasında iliřkiler gzlendi.

6. Sonu

- Veri temiz, eksiksiz ve modele hazır hle getirildi
- NLP tabanlı zellikler ile metin verisi sayısala dnřtrld
- Grselleřtirmeler ile `TedaviSuresi` iliřkili deęiřkenler belirlendi

Gzlem: Bu rapor, ileride modelleme veya istatistiksel analiz iin gl bir temel saęlar.