

Veri Madenciliđi

Giriř

İçerik

Kapanıř

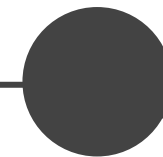
Merhaba.

Emre Can Öner

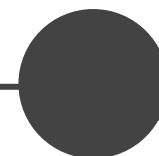
–

emrecanoner@outlook.com

Social Media Sentiment Analysis II



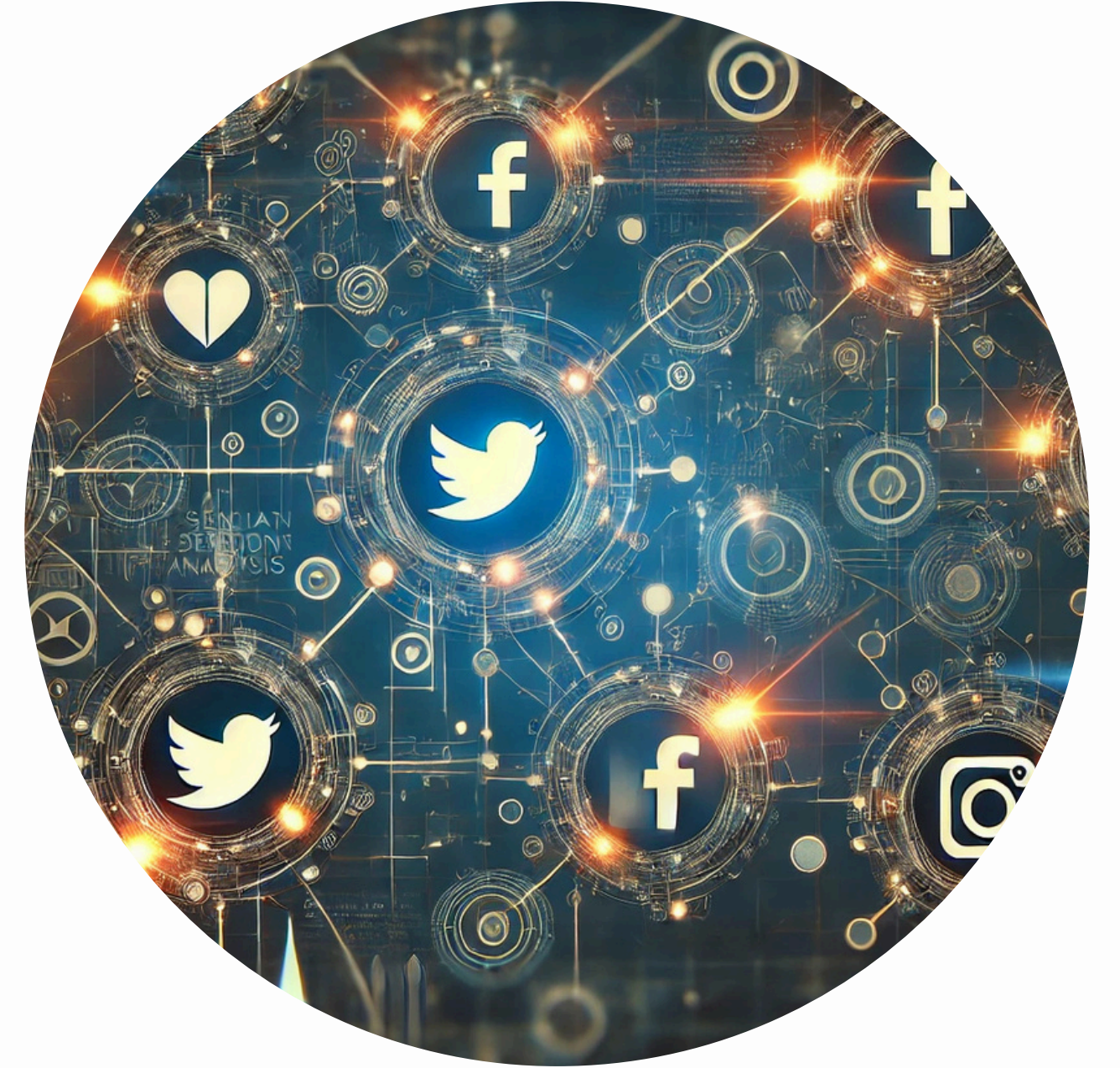
Veri Seti İşleme ve Hazırlık



Veri Seti

Orijinal Veri Seti

- Toplam örnek sayısı: 732
- Dağılım:
 - * Nötr: 614 örnek
 - * Pozitif: 98 örnek
 - * Negatif: 20 örnek
- 17 farklı özellik içeren yapılandırılmış veri



Veri Dengелеme Sureci

SMOTE Tekniđi (Synthetic Minority Over-sampling Technique)

AMA: Az temsil edilen sınıfların rneklem sayısını artırmak

NASIL?: Mevcut rneklerden yeni, sentetik rnekler oluřturur

NEDEN?: Model yanlılıđını nler ve sınıflandırma performansını iyileřtirir

- Az temsil edilen sınıflar iin sentetik rnekler
- Her sınıf iin 100 rnek hedefi
- Toplam 300 dengeli rnek

Veri Dengeleme Süreci

Veri Zenginleştirme Teknikleri

EMOJİ VARYASYONLARI

AMAÇ: Duygu ifadelerini güçlendirmek

NASIL?: Duygu sınıfına uygun emojiler eklenir

NOKTALAMA İŞARETLERİ

AMAÇ: Metin tonunu zenginleştirmek

NASIL?: Duyguya uygun noktalama işaretleri eklenir

İstatistiksel Analiz Metodolojisi

Temel Analizler

Spearman Korelasyon Analizi

AMAÇ: Değişkenler arası ilişkileri ölçmek

NEDEN SPEARMAN?:

- Normal dağılım varsayımı gerektirmez
- Sıralı veriler için uygundur

YORUMLAMA:

- -1 ile +1 arası değerler
- $p < 0.05$ anlamlılık seviyesi

Effect Size (Cohen's d)

AMAÇ: Gruplar arası farkların büyüklüğünü ölçmek

YORUMLAMA:

- 0.2: Küçük etki
- 0.5: Orta etki
- 0.8: Büyük etki
- >1.0: Çok büyük etki

ÖNEMİ: İstatistiksel anlamlılığın pratik önemini gösterir

İleri Düzey Analizler

MANOVA (Multivariate Analysis of Variance)

AMAÇ: Birden fazla bağımlı değişkenin analizi
NEDEN?: Duygu bileşenlerinin (pozitif, negatif, nötr) birlikte incelenmesi

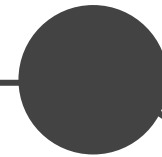
METRİKLER:

- Wilks' lambda: Gruplar arası farkların gücü
- F-değeri: Varyans analizi sonucu
- p-değeri: İstatistiksel anlamlılık

Post-hoc Tukey HSD

AMAÇ: Gruplar arası detaylı karşılaştırma
NEDEN?: Hangi grupların birbirinden farklı olduğunu belirlemek
FWER (0.05): Çoklu karşılaştırmalarda hata oranı kontrolü

Güvenilirlik ve Geçerlilik



Güvenilirlik Testleri

Cronbach's Alpha

AMAÇ: İç tutarlılık ölçümü

YORUMLAMA:

- <0.5 : Düşük güvenilirlik
- $0.5-0.7$: Orta güvenilirlik
- >0.7 : Yüksek güvenilirlik

Test-retest Güvenilirliđi

AMAÇ: Sonuçların zaman içindeki tutarlılıđı

NASIL?: Aynı örneklem üzerinde tekrar ölçüm

YORUMLAMA: >0.7 ideal deđer

İstatistiksel Güç Analizi

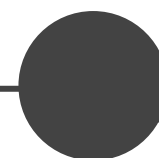
AMAÇ: Örneklem büyüklüğünün yeterliliğini değerlendirme

NASIL?: Farklı effect size değerleri için güç hesaplama

YORUMLAMA:

- >0.80: Yeterli güç
- >0.90: Mükemmel güç

Temel İstatistiksel Analiz Sonuçları



Spearman Korelasyon Analizleri

Metin Uzunluğu İlişkileri

Nötr Duygu ile: $r = 0.330$, $p < 0.0001$

- YORUM: Orta düzeyde pozitif ilişki
- ANLAMI: Uzun metinler daha çok nötr duygu içeriyor

Negatif Duygu ile: $r = -0.138$, $p = 0.0168$

- YORUM: Zayıf negatif ilişki
- ANLAMI: Kısa metinlerde negatif duygular daha yoğun

Kelime Sayısı İlişkileri

Bileşik Duygu ile: $r = 0.135$, $p = 0.0196$

- YORUM: Zayıf pozitif ilişki
- ANLAMI: Kelime sayısı arttıkça pozitif duygu eğilimi artıyor

Nötr Duygu ile: $r = 0.272$, $p < 0.0001$

- YORUM: Orta düzeyde pozitif ilişki
- ANLAMI: Çok kelimeli metinler daha nötr

Effect Size Analizi

Pozitif vs Negatif Duygular

Cohen's d = 4.307

- YORUM: Çok büyük etki büyüklüğü
- ANLAMI: İki duygu sınıfı arasında çok belirgin ayrım
- ÖNEMİ: Sınıflandırma modelinin güvenilirliğini destekler

Effect Size Analizi

Pozitif vs Nötr Duygular

Cohen's d = 1.365

- YORUM: Büyük etki büyüklüğü
- ANLAMI: Pozitif ve nötr metinler net şekilde ayrılıyor
- UYGULAMA: Sınıflandırma için güvenilir bir metrik

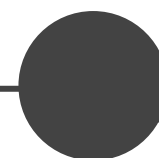
Effect Size Analizi

Negatif vs Nötr Duygular

Cohen's d = -1.452

- YORUM: Büyük negatif etki
- ANLAMI: Negatif ve nötr metinler belirgin şekilde farklı
- ÇIKARIM: Duygu sınıfları arasında güçlü ayrım

İleri Düzey Analiz Sonuçları



MANOVA Sonuçları

Genel Model Deęerlendirmesi

Wilks' lambda = 0.2853

- YORUM: Çok güçlü model uyumu
- ANLAMI: Duygu sınıfları arasında anlamlı farklılık

F-deęeri = 85.7677, $p < 0.0001$

- YORUM: Yüksek derecede anlamlı
- GÜVEN: %99.99'dan yüksek güven düzeyi

Çoklu Deęişken Etkileri

Pillai's trace = 0.8540

- YORUM: Güçlü çoklu deęişken etkisi
- ANLAMI: Duygu bileřenleri birlikte anlamlı



Post-hoc Analiz Sonuçları

Gruplar Arası Karşılaştırmalar

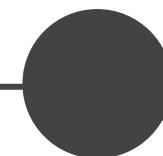
Negatif-Nötr Farkı

- Ortalama fark = 0.6719
- $p < 0.001$
- YORUM: Anlamlı pozitif fark

Negatif-Pozitif Farkı

- Ortalama fark = 1.2244
- $p < 0.001$
- YORUM: En büyük grup farkı

Güvenilirlik ve Model Performansı



Güvenilirlik Metrikleri

Cronbach's Alpha = 0.129

YORUM: Düşük iç tutarlılık

NEDEN?: Duygu ifadelerinin karmaşık yapısı

ÖNERİ: Gelecek çalışmalarda iyileştirme gerekli

Test-retest Güvenilirliği = 0.169

YORUM: Düşük tekrar tutarlılığı

ÇIKARIM: Duygu analizi sonuçları değişkenlik gösterebilir

ÇÖZÜM ÖNERİSİ: Daha güçlü modeller geliştirilebilir

İstatistiksel Güç

Güç Değerleri

Effect size 0.2 için: Power = 0.932

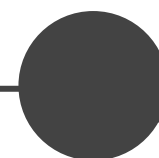
Effect size 0.5 için: Power = 1.000

Effect size 0.8 için: Power = 1.000

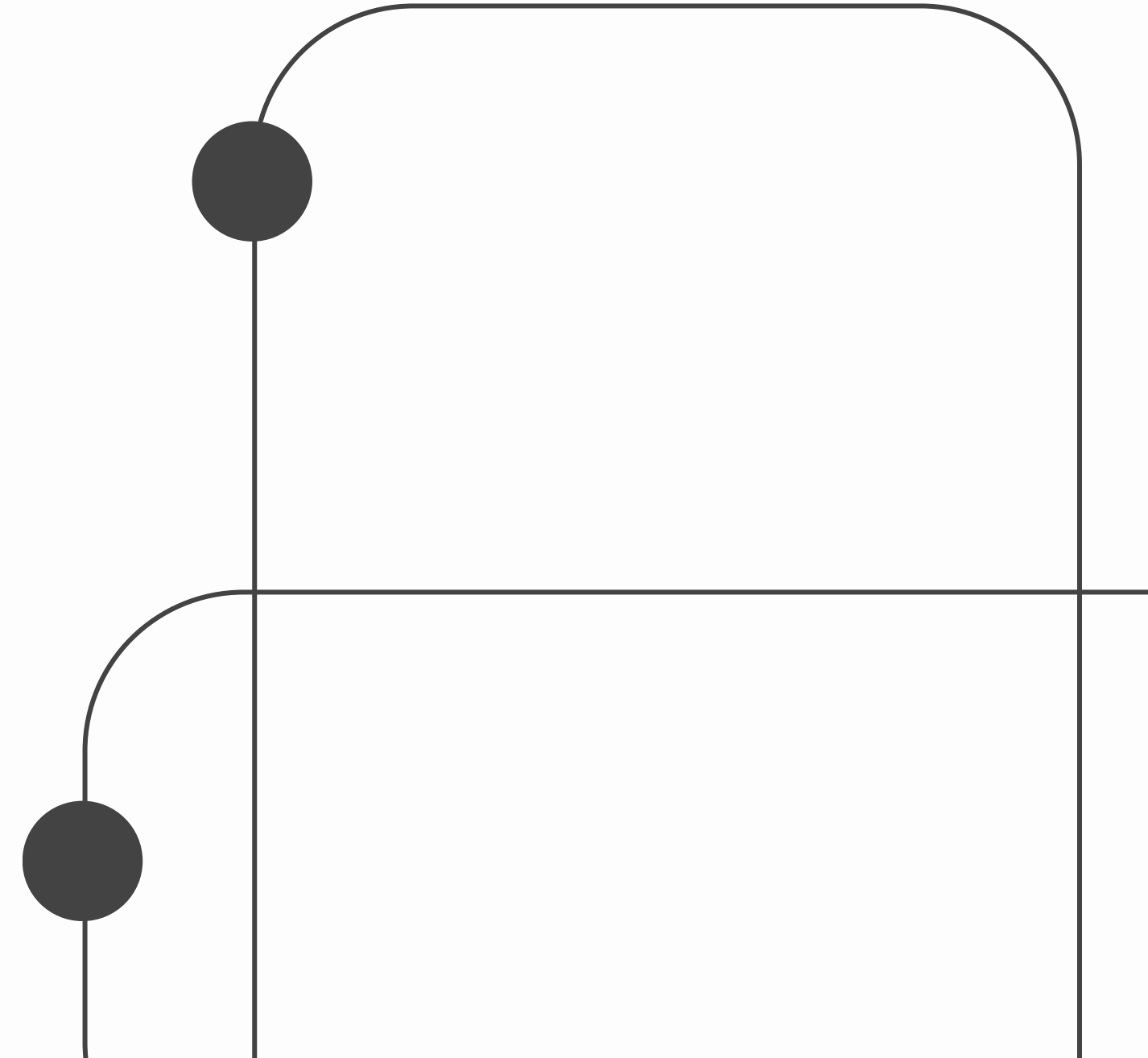
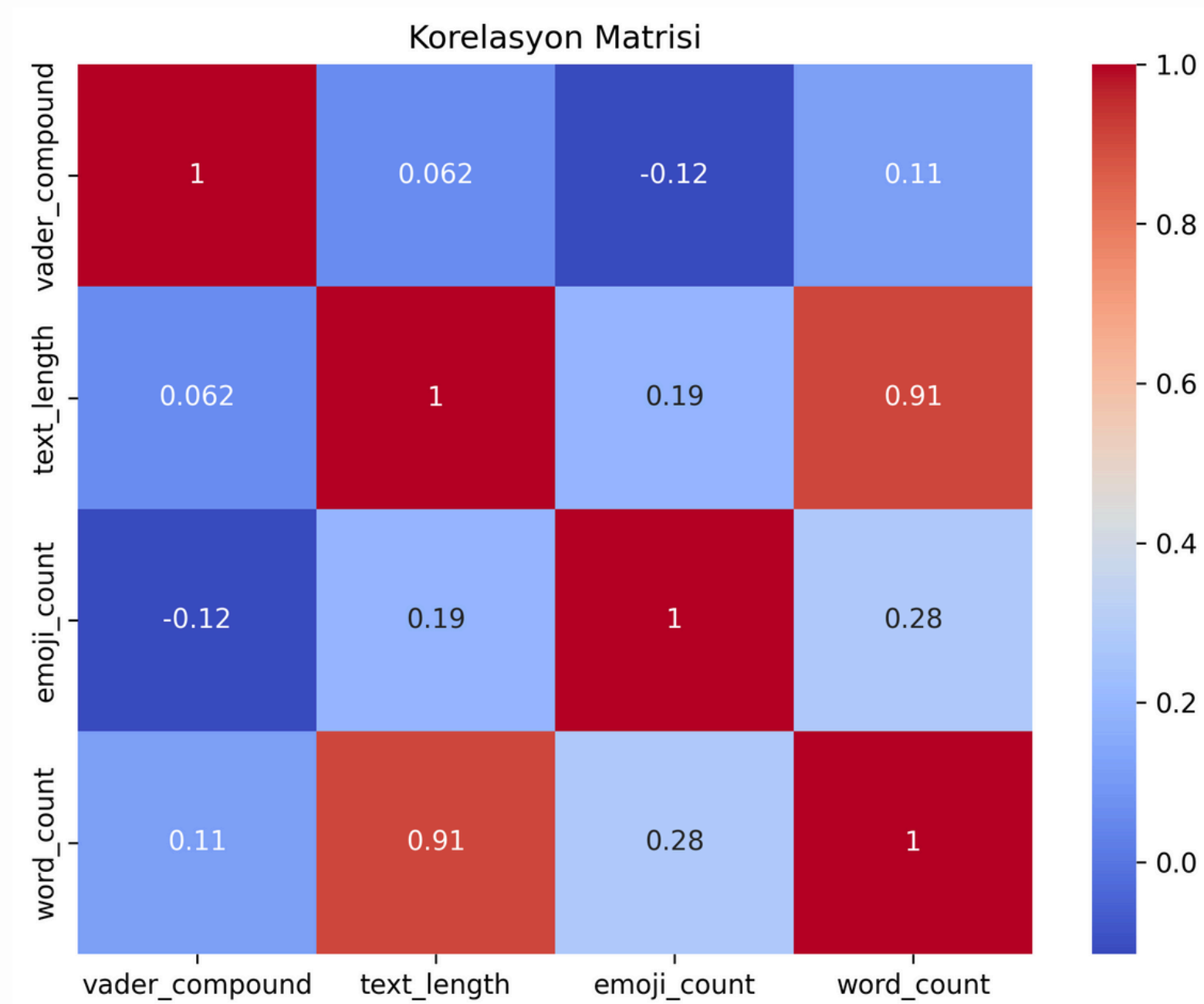
YORUM: Mükemmel istatistiksel güç

ANLAMI: Örneklem büyüklüğü yeterli

VERİ ANALİZİ GÖRSEL SONUÇLARI



Korelasyon Matrisi (Heatmap)



Korelasyon Matrisi Analizi

Metin Uzunluğu ve Kelime Sayısı İlişkisi

- Çok güçlü pozitif korelasyon ($r = 0.91$)
- Bu beklenen bir sonuç, çünkü uzun metinler doğal olarak daha çok kelime içerir
- Neredeyse mükemmel doğrusal ilişki

Duygu Skoru (VADER Compound) İlişkileri

- Metin uzunluğu ile zayıf pozitif ($r = 0.062$)
- Kelime sayısı ile zayıf pozitif ($r = 0.11$)
- Emoji sayısı ile zayıf negatif ($r = -0.12$)
- Duygu skorunun diğer özelliklerle zayıf ilişkisi, duygunun bu metriklerden bağımsız olduğunu gösterir

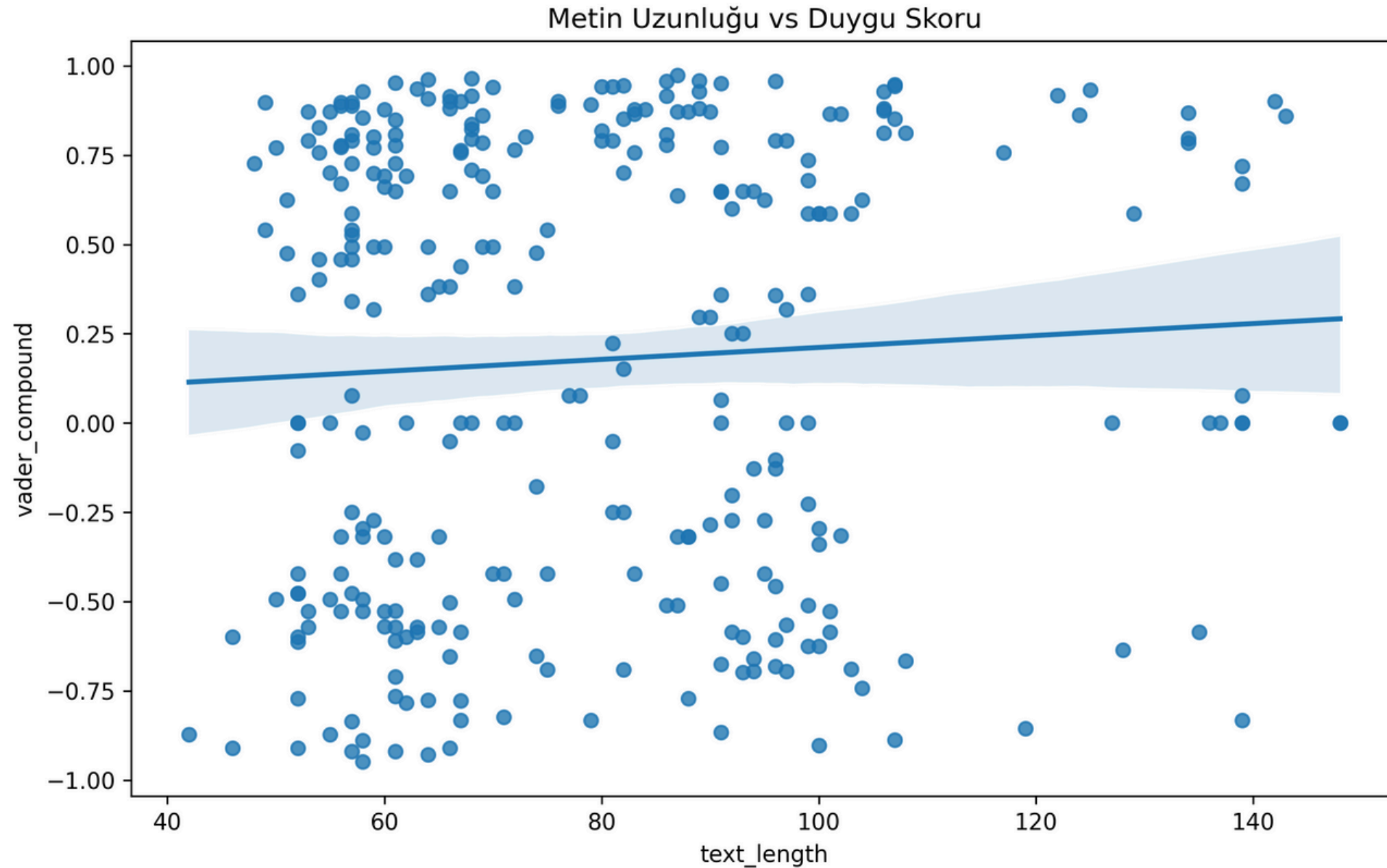
Emoji Kullanımı

- Kelime sayısı ile zayıf pozitif ilişki ($r = 0.28$)
- Metin uzunluğu ile zayıf pozitif ilişki ($r = 0.19$)
- Emoji kullanımının metin uzunluğundan çok etkilenmediğini gösterir

Pratik Çıkarımlar:

- Duygu analizi yaparken metin uzunluğu normalize edilmeli
- Emoji sayısı bağımsız bir faktör olarak değerlendirilmeli
- Kelime sayısı ve metin uzunluğundan sadece birini modele dahil etmek yeterli olabilir

Metin Uzunluğu ve Duygu İlişkisi (Regression Plot)



Regresyon Analizi

Genel Trend

- Çok hafif pozitif eğim (neredeyse yatay)
- Mavi çizgi regresyon doğrusunu gösteriyor
- Açık mavi alan güven aralığını temsil ediyor

Dağılım Özellikleri

- Noktalar -1 ile +1 arasında geniş bir dağılım gösteriyor
- Yoğunlaşma özellikle -0.5 ile +0.75 aralığında
- Metin uzunluğu genellikle 40-140 karakter arasında

İlişki Analizi

- Zayıf pozitif korelasyon
- Yüksek varyans (noktaların regresyon çizgisinden uzak dağılımı)
- Güven aralığının genişliği tahmin belirsizliğini gösteriyor

Önemli Gözlemler

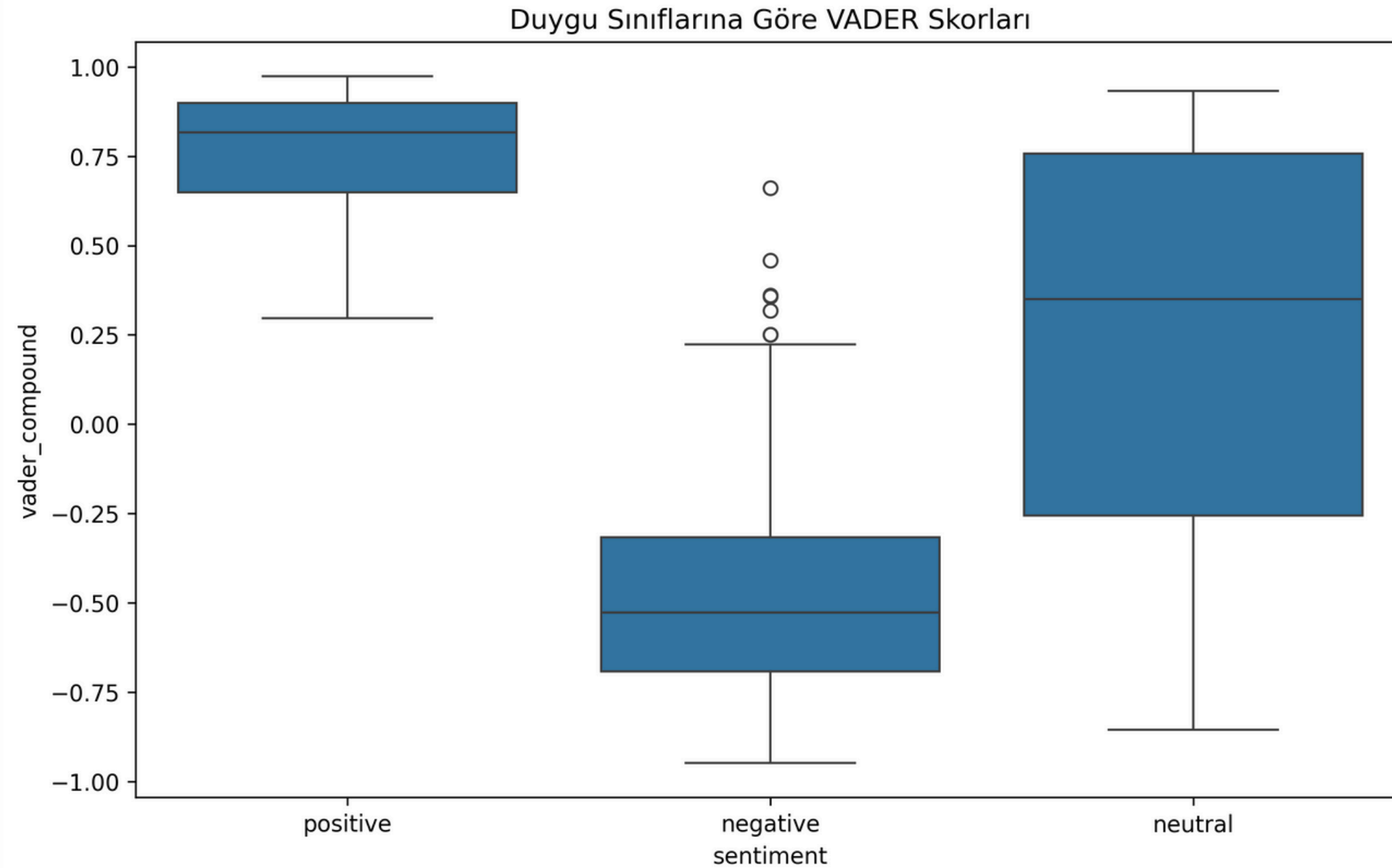
- Kısa metinlerde (-40-60 karakter) duygu skorları daha değişken
- Uzun metinlerde (100+ karakter) dağılım biraz daha az yoğun
- Uç değerler her metin uzunluğunda mevcut

Pratik Çıkarımlar:

- Metin uzunluğu duygu skorunu önemli ölçüde etkilemiyor
- Duygu analizi için metin uzunluğu normalize edilmeli
- Kısa metinlerde daha dikkatli analiz gerekebilir
- Model geliştirirken metin uzunluğu tek başına yeterli bir faktör değil

Bu analiz, korelasyon matrisinde gördüğümüz zayıf ilişkiyi ($r = 0.062$) görsel olarak da doğruluyor.

Duygu Sınıflarına Göre VADER Skorları (Box Plot)



Duygu Sınıflarına Göre VADER Skorları Analizi

Pozitif Sınıf Özellikleri

- Medyan: ~0.80 civarında
- Dar kutu aralığı: Tutarlı pozitif skorlar
- Alt sınır: ~0.30
- Üst sınır: ~0.95
- Az sayıda aykırı değer

Negatif Sınıf Özellikleri

- Medyan: ~-0.50 civarında
- Orta genişlikte kutu: Makul tutarlılık
- Birkaç pozitif aykırı değer (0.25-0.65 aralığında)
- Alt sınır: ~-0.95
- Üst sınır: ~-0.25

Nötr Sınıf Özellikleri

- Medyan: ~0.35 civarında
- Geniş kutu aralığı: Yüksek varyans
- Geniş dağılım: -0.85 ile 0.95 arası
- En fazla değişkenlik gösteren sınıf

Önemli Gözlemler:

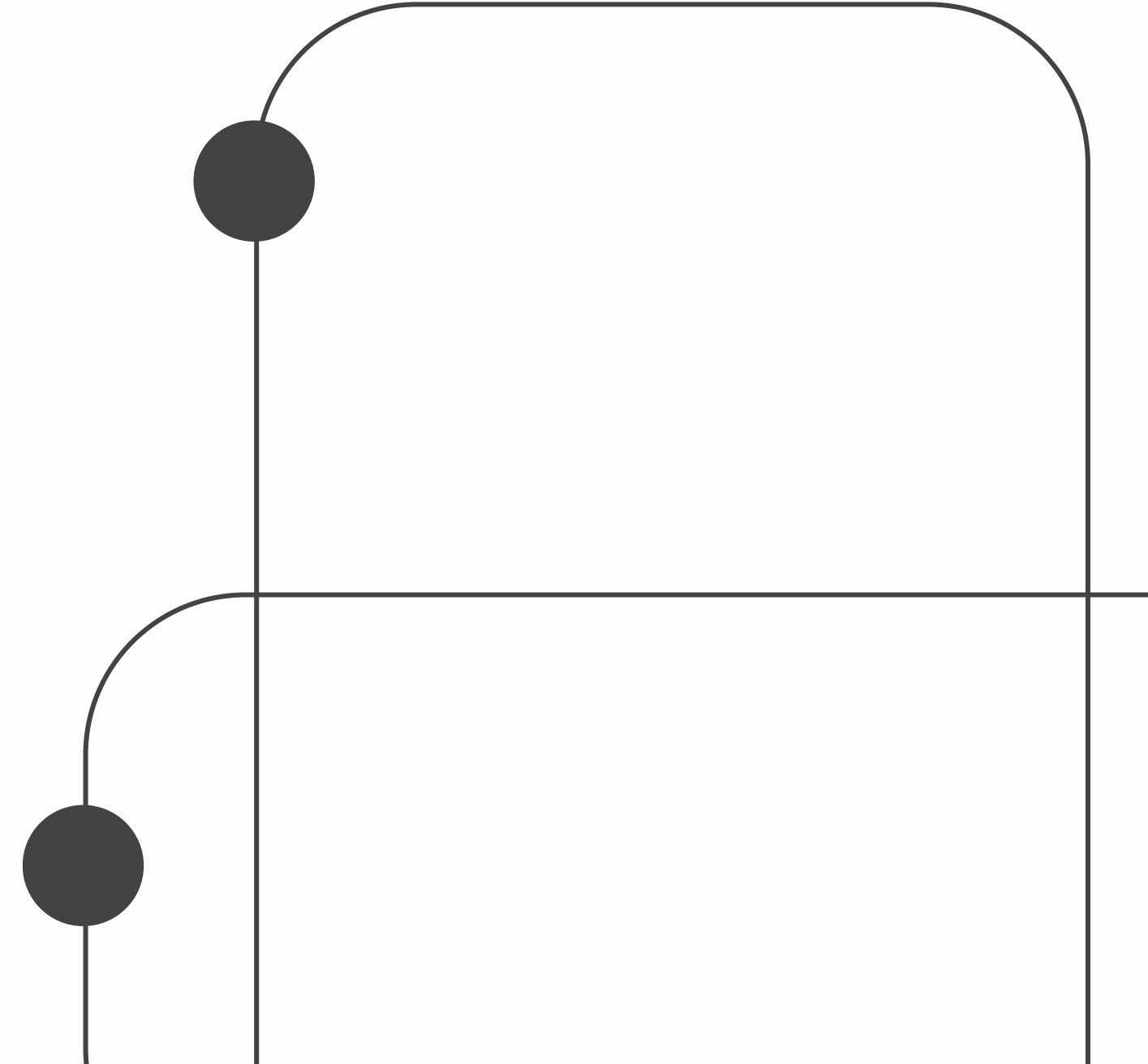
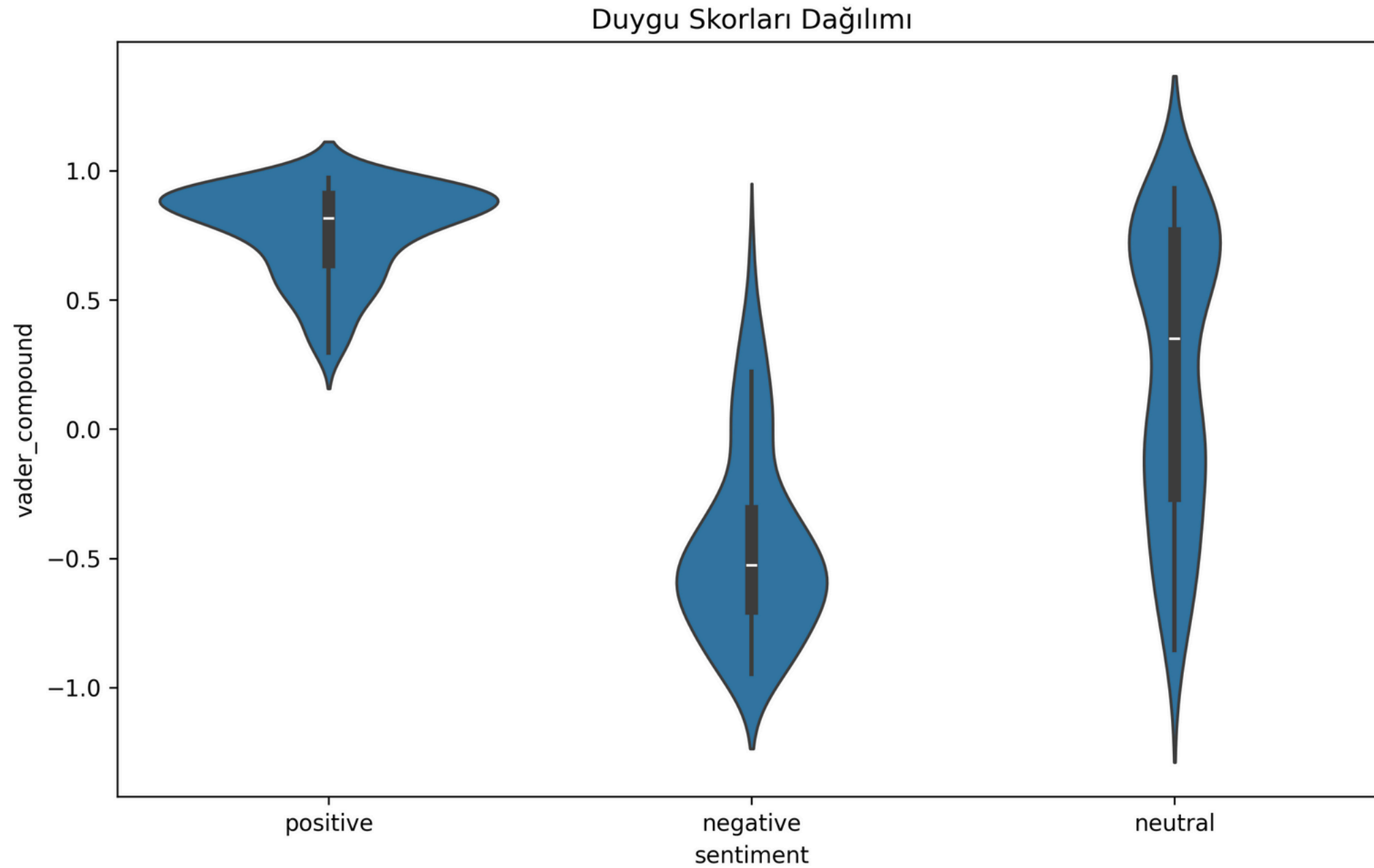
- Sınıflar arası belirgin ayırım var
- Pozitif sınıf en tutarlı sonuçları veriyor
- Nötr sınıf en geniş dağılıma sahip
- Negatif sınıfta pozitif yönde aykırı değerler mevcut

Pratik Çıkarımlar:

- VADER pozitif duyguları daha iyi tespit ediyor
- Nötr metinlerin sınıflandırılması daha zor
- Negatif metinlerde yanlış pozitif olasılığı var
- Sınıflandırma için eşik değerleri dikkatli seçilmeli

Bu analiz, duygu sınıflandırma modelinin performansını ve potansiyel iyileştirme alanlarını net bir şekilde gösteriyor.

Duygu Skorları Dağılımı (Violin Plot)



Duygu Skorları Dağılımı Analizi

Pozitif Sınıf Dağılımı

- Çan şeklinde yoğunlaşma (~0.75 civarında)
- Alt kısımda incelme (0.2-0.4 arası)
- Üst kısımda geniş dağılım (0.8-1.0 arası)
- Simetrik olmayan dağılım
- Medyan çizgisi ~0.8 civarında

Negatif Sınıf Dağılımı

- Çift tepeli dağılım
- Ana yoğunluk -0.5 ile -0.8 arasında
- İkincil yoğunluk -0.2 civarında
- En dar dağılıma sahip sınıf
- Medyan çizgisi ~-0.5 civarında

Nötr Sınıf Dağılımı

- Uzun ve ince dağılım
- Geniş bir aralığa yayılım (-1.0 ile 1.0 arası)
- Orta kısımda hafif şişkinlik
- Medyan çizgisi ~0.35 civarında
- En geniş dağılıma sahip sınıf

Önemli Gözlemler:

- Her sınıfın kendine özgü bir dağılım şekli var
- Pozitif sınıf en belirgin ve tutarlı dağılımı gösteriyor
- Negatif sınıf iki farklı yoğunluk bölgesi içeriyor
- Nötr sınıf en belirsiz ve geniş dağılıma sahip

Pratik Çıkarımlar:

- Pozitif duygu tespiti daha güvenilir
- Negatif duygu tespitinde alt kategoriler olabilir
- Nötr sınıflandırma için daha hassas kriterler gerekli
- Sınıflandırma eşikleri dağılım şekillerine göre ayarlanmalı

Bu violin plot, box plot'ta gördüğümüz bulguları daha detaylı bir şekilde doğruluyor ve dağılımların şekillerini daha net gösteriyor.

Veri Madenciliđi

Giriř

İçerik

Kapanıř

Teřekkürler.

Emre Can Öner

–

emrecanoner@outlook.com