

Konvolüsyonel Sinir Ağları ile Konuşma Duygusu Tanıma Speech Emotion Recognition with Convolutional Neural Networks

Ömer Faruk ORUÇ^{1,a}

¹Elektrik-Elektronik Mühendisliği, Mühendislik Fakültesi, Sakarya Üniversitesi, Serdivan, Sakarya, Türkiye

^a omer.oruc2@ogr.sakarya.edu.tr,

Özetçe—Bu çalışma, konuşma duygu tanıma (Speech Emotion Recognition, SER) alanında makinelerin insan duygularını anlaması ve bu duygulara tepki verebilmesi amacıyla Konvolüsyonel Sinir Ağları (CNN) tabanlı bir model geliştirmeyi amaçlamaktadır. Toronto Duygusal Konuşma Seti (TESS) ve diğer kapsamlı veri setleri kullanılarak konuşma sinyalleri analiz edilmiş, veri artırma teknikleriyle çeşitlendirilmiş ve Mel-frekans kepstrum katsayıları (MFCC) gibi akustik özellikler çıkarılmıştır. Model, veri işleme ve öğrenme sürecinde gürültü azaltma, zaman kaydırma ve ton değiştirme teknikleriyle desteklenmiştir. Eğitim aşamasında %99,89'a, doğrulama aşamasında ise %93,01 doğruluğa ulaşan CNN modeli, duygusal ifadeleri başarıyla sınıflandırmış ve düşük kayıp oranıyla yüksek genelleme kapasitesi göstermiştir. Bu sonuçlar, konuşma sinyallerinden duygu çıkarımında CNN tabanlı yöntemlerin etkinliğini ortaya koymaktadır. Çalışma, SER sistemlerinin müşteri hizmetleri, sanal asistanlar ve insan-makine etkileşimi alanlarında daha kişiselleştirilmiş ve anlamlı deneyimler sunma potansiyeline işaret etmektedir. Model, bu alanda gelecekteki uygulamalar için güçlü bir temel oluştururken, yapay zeka destekli duygu tanıma teknolojilerinin geliştirilmesine yönelik önemli bir katkı sunmaktadır.

Anahtar Kelimeler—konvolüsyonel sinir ağları; duygu sınıflandırması; ses tabanlı duygu tanıma; ses sinyali analizi; özellik çıkarımı

Abstract— This study aims to develop a Convolutional Neural Network (CNN)-based model in the field of Speech Emotion Recognition (SER) to enable machines to understand and respond to human emotions. Speech signals were analyzed using the Toronto Emotional Speech Set (TESS) and other comprehensive datasets, diversified through data augmentation techniques, and acoustic features such as Mel-frequency cepstral coefficients (MFCC) were extracted. The model was supported by techniques like noise reduction, time shifting, and pitch alteration during data processing and training. The CNN model achieved 99.89% accuracy in the training phase and 93.01% accuracy in the validation phase, successfully classifying emotional expressions and demonstrating high generalization capacity with a low loss rate. These results highlight the effectiveness of CNN-based methods in extracting emotions from speech signals. The study points to the potential of SER systems to provide more personalized and meaningful experiences in areas such as customer service, virtual assistants, and human-machine interaction. The model establishes a strong foundation for future applications in this field and makes a significant contribution to

the development of artificial intelligence-supported emotion recognition technologies.

Keywords—convolutional neural networks; emotion classification; speech based emotion recognition; speech signal analysis; feature extraction

I. GİRİŞ

Makinelerin insan duygularını tanıma ve bu duygulara tepki verme yeteneği, akıllı sistemlerin geliştirilmesinde giderek daha önemli hale gelmiştir [1]. Konuşma Duygu Tanıma (Speech Emotion Recognition- SER), insanlar ve makineler arasında doğal ve sorunsuz etkileşimler sağlama açısından kilit bir rol oynamakta ve sanal asistanlar, müşteri hizmetleri ve insan-bilgisayar etkileşimi gibi alanlarda büyük önem taşımaktadır [2]. Konuşma yoluyla iletilen duyguları anlayarak, akıllı sistemler daha kişiselleştirilmiş ve anlamlı tepkiler sağlayabilir, bu da onların kullanılabilirliğini ve kullanıcı deneyimini artırır.

Bu projede, SER'yi geliştirmek için Konvolüsyonel Sinir Ağı (CNN) tabanlı bir yaklaşım uygulandı. Model, 20 epoch boyunca eğitilmiş olup, başlangıç epoch'unda %55,28 olan eğitim doğruluğu, son epoch'ta %99,89'a kadar yükselmiştir. Doğrulama doğruluğu %93,01'e ulaşarak modelin yeni verilere iyi bir şekilde genelleme yapabildiğini göstermiştir. Ayrıca, doğrulama kaybı düzenli bir şekilde azalarak 0,2810 seviyesine inmiş ve bu durum modelin güçlü performansını ve kararlılığını ortaya koymuştur. Eğitim sırasında öğrenme oranının azaltılması gibi teknikler, daha hızlı yakınsama ve daha iyi optimizasyon sağlamak amacıyla kullanılmıştır.

Bu sonuçlar, CNN tabanlı yaklaşımın konuşma verilerinden duyguları tanımadaki etkinliğini göstermektedir. Bulgular, SER uygulamalarında daha fazla ilerlemeye olanak tanıyarak, insan duygularını daha iyi anlayabilen ve bu duygulara uygun tepkiler verebilen duygusal olarak farkındalıklı teknolojilerin geliştirilmesine katkıda bulunmaktadır.

II. MATERYAL VE METOT

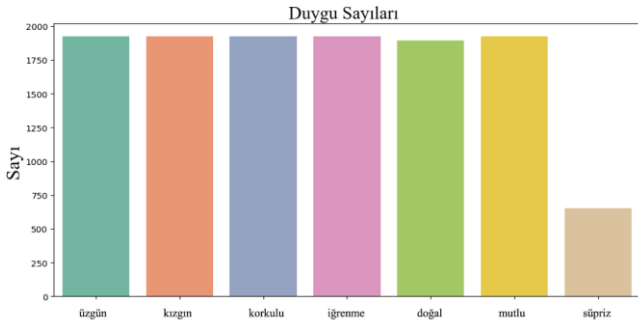
A. Veriseti

Veri seti olarak Toronto Duygusal Konuşma Seti (TESS), Toronto Üniversitesi Psikoloji Bölümü tarafından duygu tanıma araştırmaları için geliştirilen açık kaynaklı bir veri seti kullanılmıştır. Bu veri seti, 390 kbps bit hızında kaydedilmiş toplamda 2.800 yüksek kaliteli .wav ses dosyasından oluşmaktadır. Veri kayıtları, 26 ve 64 yaşlarında olan, ana dili

İngilizce olan iki kadın seslendirme sanatçısı tarafından gerçekleştirilmiştir. Sanatçılar, "speak the word" taşıyıcı ifadesini söyledikten sonra 200 hedef tek kelimedenden birini dile getirmiştir. Veri seti; iğrenme, öfke, mutluluk, korku, nötrlük, üzüntü ve hoş bir sürpriz olmak üzere yedi farklı duygu kategorisini kapsamaktadır.

TESS veri seti, dört ana klasörde düzenlenmiş daha geniş bir ses duygu veri seti gurubunun alt parçalarıdır ve bunlar CREMA-D (/content/Crema), RAVDESS (/content/Ravdess), SAVEE (/content/Savee) ve TESS (/content/Tess) dir. TESS klasörü içinde ses dosyaları, duygu kategorilerine göre gruplandırılmıştır ve her kategori farklı ses dosyaları içermektedir. Dosyalar, belirli bir adlandırma kuralına göre etiketlenmiştir: "OAF" ile başlayan dosya adları yaşlı sanatçının kayıtlarını, "YAF" ile başlayanlar ise genç sanatçının kayıtlarını ifade etmektedir. Her dosya adı, bir sanatçı kodu, hedef kelime ve ifade edilen duyguyu içermektedir.

TESS, konuşmada duygu tanıma, ses tabanlı duygu algılama için makine öğrenimi modellerinin eğitimi, insan-bilgisayar etkileşimlerinin geliştirilmesi, konuşma sentezi ve modülasyon araştırmaları gibi uygulamalarda yaygın olarak kullanılmaktadır. TESS, CREMA-D, RAVDESS ve SAVEE gibi diğer veri setleriyle birleştirilerek araştırmacılara çeşitli ses örneklerine erişim imkanı sunmakta, bu da çeşitli ses duygu tanıma görevleri için sağlam ve genelleştirilebilir modellerin geliştirilmesini sağlamaktadır. Verisetinin görsel olarak yedi duygu kategorisindeki veri dağılımı Şekil I'de gösterilmektedir.



Şekil I. Veriseti dağılımı

B. Veri Ön İşleme

Yapay zeka modellerinin başarılı bir şekilde eğitilebilmesi için kullanılan ses verilerinin doğru şekilde ön işlenmesi kritik bir adımdır [3]. Bu nedenle, ses verileri üzerinde ilk olarak çevresel gürültülerin ve istenmeyen seslerin etkisini minimize etmek için bir dizi işlem gerçekleştirilmiştir. İnsan işitme aralığında olmayan frekanslar (16 Hertz ile 16.000 Hertz dışında kalan sesler), band geçiren filtreler kullanılarak temizlenmiştir. Bu işlem hem verinin kalitesini artırmak hem de modelin yalnızca gerekli bilgilere odaklanmasını sağlamak için gerçekleştirilmiştir. Gürültülerin etkisini ortadan kaldırarak, modelin daha doğru ve güvenilir sonuçlar üretmesi hedeflenmiştir.

Verilerin hazırlanması aşamasında, CREMA, RAVDESS, SAVEE ve TESS veri setlerinden faydalanılmıştır. Her bir veri seti, detaylı olarak analiz edilerek dosya yolları ve duygu etiketleri çıkarılmıştır. Bu süreçte, verilerin sistematik bir

şekilde ayrıştırılması ve düzenlenmesi, modelin eğitimi için gerekli temel yapıyı oluşturmuştur.

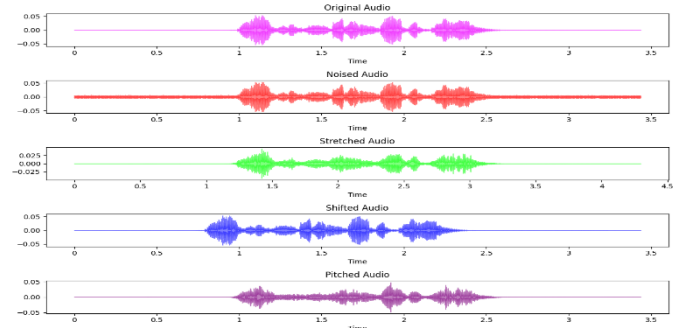
Modelin dayanıklılığını ve esnekliğini artırmak için veri artırma yöntemleri yoğun bir şekilde uygulanmıştır. Bu yöntemlerin temel amacı, modelin farklı veri türleri ve durumlarına uyum sağlayabilmesini sağlamaktır. Gürültü ekleme işlemi, gerçek dünyadaki ses verilerinde karşılaşılabilecek çeşitli çevresel gürültüleri simüle etmek için kullanılmıştır. Gaussian tipi rastgele gürültü, ses sinyallerine belirli bir eşik veya rastgele oranlarda eklenmiş, böylece modelin bu tür etkilerle başa çıkma kabiliyeti geliştirilmiştir.

Zaman ekseninde yapılan kaydırma işlemi (shifting), ses sinyallerinin başlangıç noktalarında değişiklikler yaparak veri setine çeşitlilik kazandırmıştır. Bu yöntem, özellikle sinyallerin zamanlamasında meydana gelebilecek farklılıkları simüle ederek, modelin bu tür varyasyonlara karşı daha duyarlı hale gelmesini sağlamaktadır [4]. Perde değiştirme (pitching) tekniği ile sesin tonunda değişiklikler yapılmış, belirli bir aralıkta rastgele seçilen bir faktörle ton artırılmış ya da düşürülmüştür. Böylece, modelin farklı konuşma tarzlarını ve duygusal tonlamaları daha kolay ayırt edebilmesi sağlanmış oldu.

Ses sinyallerine uygulanan zaman uzatma (stretching) işlemi ise, konuşma hızındaki değişimleri yansıtarak veri setine önemli bir çeşitlilik eklemiştir. Bu işlemde, ses sinyalleri ya uzatılmış ya da sıkıştırılmış, bu sayede modelin hızlı veya yavaş konuşmaları daha doğru bir şekilde işlemesi sağlanmıştır. Özellikle konuşma hızına dayalı varyasyonların modellenmesi, bu işlemin model performansı üzerindeki etkisini artırmıştır.

Tüm bu veri artırma yöntemleri, yalnızca modelin genelleme kapasitesini geliştirmekle kalmamış, aynı zamanda daha sağlam ve gerçek dünya koşullarına uygun bir yapıya kavuşmasını sağlamıştır. Gürültülerin eklenmesi, sinyallerin kaydırılması, tonlamaların değiştirilmesi ve zaman uzatılması gibi işlemler, modelin yalnızca akademik bir çerçevede değil, pratik uygulamalarda da başarılı olmasına katkı sunmuştur.

Şekil II'de gösterilen görsel, uygulanan veri artırma yöntemlerinin ses sinyalleri üzerindeki etkisini net bir şekilde ortaya koymaktadır. Her bir satır, farklı bir veri artırma tekniğini temsil etmekte ve bu işlemlerin modelin eğitimi için nasıl bir çeşitlilik sağladığını göstermektedir.



Şekil II. Ses verisi ön işleme adımları

C. Özellik Çıkarımı

Ses sinyallerinin işlenmesi sırasında, duygusal durumların sınıflandırılması için akustik, istatistiksel ve zaman serisi özellikleri çıkarılmıştır. Bu süreçte, temel akustik özellikler

arasında Mel-frekans kepstrum katsayıları (MFCC), sıfır geçiş oranı (Zero Crossing Rate, ZCR), kök ortalama kare enerji (Root Mean Square Energy, RMSE) ve spektral özellikler yer almıştır. Ayrıca, zaman serisi analizi ile sinyallerin zamana bağlı değişimleri incelenmiş, bu değişimlerin frekans içeriği üzerindeki etkileri analiz edilmiştir.

İstatistiksel olarak, sinyallerin merkezi eğilimleri ve dağılım özellikleri hesaplanmıştır. Ortalama, standart sapma, maksimum ve minimum değerler gibi temel istatistiksel ölçütler öznitelik vektörlerine dahil edilmiştir. Öte yandan, sinyal işleme metrikleri kapsamında, darbe faktörü, tepe faktörü ve sinyal-gürültü oranı (SNR) gibi metrikler de hesaplanarak, veri setinin daha ayrıntılı bir temsili sağlanmıştır. Bu özelliklerin çıkarılması, sinyallerin anlamlı bir şekilde temsil edilmesine ve sınıflandırma modellerinin etkinliğinin artırılmasına olanak sağlamıştır.

D. Sınıflandırma Modelinin Geliştirilmesi

Sınıflandırma modeli, derin öğrenme yöntemlerinden biri olan Konvolüsyonel Sinir Ağı (CNN) ile oluşturulmuştur. CNN mimarisi, giriş verilerinden karmaşık özellikler öğrenme ve çıkarma yeteneği ile dikkat çekmektedir. Model, sıralı bir yapı izleyerek, 512 filtreye sahip bir 1D konvolüsyon katmanı ile başlamış, ardından küme normalizasyonu (batch normalization) ve maksimum havuzlama (max-pooling) işlemleri uygulanmıştır. Daha derin katmanlarda filtre sayısı kademeli olarak azaltılmış ve her katmanda ReLU aktivasyon fonksiyonu kullanılmıştır.

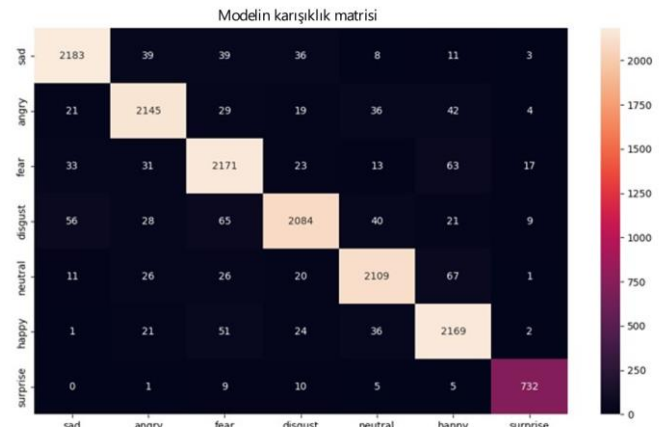
Modelin son aşamasında, düzleştirilmiş (flattening) özellik haritaları tam bağlantılı katmanlara geçirilmiş ve softmax aktivasyon fonksiyonu ile sınıflandırma gerçekleştirilmiştir. Modelin eğitimi sırasında doğruluk oranı istikrarlı bir şekilde artarak 20 epoch sonunda %93 doğrulama doğruluğuna ulaşmıştır.

E. Performans Değerlendirme Kriterleri

Bu çalışmada, konuşma sinyallerinden duyguları tanımlamak amacıyla kapsamlı bir ses sinyali analizi gerçekleştirilmiştir. Analiz süreci, veri setinin pandas kütüphanesi kullanılarak yüklenmesi ve işlenmesiyle başlamıştır. Eksik değerler, analizin doğruluğunu korumak için sıfırlarla doldurulmuş ve veri setinin yapısı incelenmiştir. İnceleme sonucunda, her bir örneğin 2.377 özellik içerdiği ve toplamda 48.648 örnekten oluştuğu belirlenmiştir.

Ses sinyallerinden Mel-frekans kepstrum katsayıları (MFCC), sıfır geçiş oranı (zero crossing rate) ve RMS enerji gibi temel akustik özellikler çıkarılmış ve bu özellikler bir CSV dosyasında saklanmıştır.

Daha sonra, veri seti giriş özellikleri (X) ve hedef etiketler (Y) olmak üzere iki ayrı bileşene ayrılmıştır. Duygusal durumları temsil eden hedef etiketler, LabelEncoder ile sayısal formatta kodlanmış ve ardından Keras kütüphanesinin to_categorical fonksiyonu kullanılarak kategorik değerlere dönüştürülmüştür. Bu işlem sonucunda, sınıflandırma için yedi farklı duygu kategorisi tanımlanmıştır: 'öfkeli' (angry), 'iğrenme' (disgust), 'korku' (fear), 'mutlu' (happy), 'nötr' (neutral), 'üzgün' (sad) ve 'sürpriz' (surprise).



Şekil III. Karışıklık matrisi

Şekil III'de matrisin diyagonal elemanları, doğru tahmin edilen örnek sayısını göstermekte olup modelin sınıflandırma performansını temsil eder. Diyagonal dışındaki elemanlar ise yanlış sınıflandırmaları ifade etmektedir.

III. SONUÇLAR

Bu çalışmada kullanılan Konvolüsyonel Sinir Ağı (CNN, Convolutional Neural Network) mimarisi, giriş ses verilerinden anlamlı özellikler çıkartmak amacıyla çok sayıda konvolüsyon katmanı ile tasarlanmıştır. Model, sıralı bir yapı izleyerek, 512 filtre ve belirli bir çekirdek boyutuna (kernel size) sahip 1D konvolüsyon katmanı ile başlamaktadır. Bu katman, giriş verilerindeki temel özellikleri öğrenmek için tasarlanmış olup, ardından küme normalizasyonu (batch normalization) ile verilerin normalize edilmesi ve maksimum havuzlama (max-pooling) işlemiyle önemli özelliklerin öne çıkarılması sağlanmıştır.

Daha sonraki katmanlarda, modelin farklı soyutlama düzeylerinden öğrenmesini sağlamak için filtre sayısı kademeli olarak 512'den 256'ya ve ardından 128'e düşürülmüştür. Her konvolüsyon katmanında, modele doğrusal olmayanlık kazandırmak ve daha karmaşık desenleri öğrenmek amacıyla ReLU (Rectified Linear Unit) aktivasyon fonksiyonu kullanılmıştır [5]. Maksimum havuzlama katmanları, 3 veya 5 boyutlu havuzlama pencerelerine (pooling window) sahiptir ve bu pencereler, uzamsal boyutları azaltarak hesaplama yükünü düşürmenin yanı sıra aşırı öğrenmeyi (overfitting) önlemeye katkıda bulunmaktadır [6].

Modelin son aşamasında, konvolüsyon katmanlarından elde edilen özellik haritaları düzleştirilmiş (flattening) ve tam bağlantılı katmanlar (fully connected layers) aracılığıyla işlenmiştir. Softmax aktivasyon fonksiyonu ile gerçekleştirilen son sınıflandırma işlemi sonucunda model, yedi duygu kategorisine ait olasılık değerlerini üretmiştir. Bu yapı, hem giriş verilerindeki temel özellikleri etkin bir şekilde yakalamak hem de modelin genel performansını artırmak amacıyla optimize edilmiştir.

Modelin doğruluğu (accuracy), doğru tahmin edilen duyguların toplam tahminlere oranı olarak değerlendirilmiştir. Eğitim sürecinin başlangıcında %55 olan doğruluk oranı, 20. epoch sonunda %99 seviyesine ulaşmıştır. Benzer şekilde, doğrulama doğruluğu (validation accuracy) da eğitim sürecinde

iyileşme göstermiş ve %93 seviyesine kadar yükselmiştir. Bu sonuçlar, modelin eğitim verilerinden etkili bir şekilde öğrenme ve yeni verilere genelleme yapma kapasitesini göstermektedir. Özellikle mutlu, üzgün ve sürpriz gibi duyguların yüksek doğruluk oranlarıyla doğru bir şekilde sınıflandırılması, modelin farklı duygusal durumları ayırt etmedeki başarısını kanıtlamaktadır.

Modelin performansı, tahmin edilen değerlerle gerçek etiketler arasındaki uyumsuzluğu ölçen kayıp fonksiyonu üzerinden değerlendirilmiştir. Eğitim kaybı sürecin başlangıcında yüksek olmasına rağmen düzenli bir azalma göstererek minimum seviyeye ulaşmıştır. Bu durum, modelin eğitim verisindeki örüntüleri öğrenme sürecini başarıyla tamamladığını ve tahmin doğruluğunu artırdığını göstermektedir.

Doğrulama kaybı ise dalgalı bir seyir izlese de genel olarak azalma eğilimi göstermiştir. Öğrenme oranının ReduceLRonPlateau stratejisi ile dinamik olarak ayarlanması, kayıp değerlerinin daha tutarlı bir şekilde azalmasına katkıda bulunmuştur [7]. Eğitim sonunda doğrulama kaybının düşük seviyelere ulaşması, modelin genelleme kapasitesinin güçlendiğini ve overfitting'in büyük ölçüde önlendiğini ifade etmektedir.

IV. TARTIŞMA VE SONUÇLAR

Konuşma tabanlı duygu tanıma, insan-bilgisayar etkileşimi, sosyal robotik ve müşteri hizmetleri gibi birçok alanda önemli bir role sahiptir. Ses verileri, konuşmacının duygusal durumunu anlamak ve bu duruma uygun tepkiler geliştirmek için etkili bir kaynak sağlamaktadır. Bu nedenle, ses tabanlı duygu tanıma sistemlerinin geliştirilmesi giderek daha fazla önem kazanmaktadır.

Bu çalışmada, konuşma verilerinden duyguların sınıflandırılması için Konvölüsyonel Sinir Ağları (CNN) kullanılmış ve %93,01 doğrulama doğruluğu elde edilmiştir. ReduceLRonPlateau stratejisi ile öğrenme oranı dinamik olarak ayarlanmış, modelin kararlı ve başarılı bir öğrenme süreci geçirmesi sağlanmıştır. Modelin yüksek doğruluk oranları ve düşük kayıp değerleri, CNN tabanlı yaklaşımın etkinliğini ortaya koymaktadır.

Literatürle yapılan karşılaştırmalar, çalışmanın CNN tabanlı mimarisi ve veri artırma teknikleriyle duygu tanıma sistemlerinde önemli bir başarı sağladığını göstermektedir [8]. Özellikle veri çeşitliliği ve artırma yöntemleri sayesinde modelin genelleme kapasitesinin artmış olduğu söylenebilir. Ancak, bazı duygu kategorilerinde karşılaşılan sınıflandırma zorlukları, literatürdeki benzer çalışmalarla paralellik göstermekte ve bu kategoriler için daha fazla farklı özellik çıkarım tekniklerinin kullanılmasını gerektirmektedir [9].

Sonuç olarak, bu çalışma, CNN tabanlı yaklaşımların duygu sınıflandırmada etkili olduğunu göstermiştir. Elde edilen bulgular, modelin eğitim sürecinde kullanılan tekniklerin başarısını ortaya koymaktadır. Gelecekteki çalışmalar, daha geniş veri setleri ve farklı kültürel dil özellikleriyle bu alandaki ilerlemeyi destekleyebilir.

KAYNAKLAR

- [1] Kusal, S., Patil, S., Kotecha, K., Aluvalu, R., and Varadarajan, V., ``AI-Based Emotion Detection for Textual Big Data: Techniques and Contribution'', Big Data and Cognitive Computing, 5(3), 2021.
- [2] Davila-Gonzalez, S., and Martin, S., ``Human Digital Twin in Industry 5.0: A Holistic Approach to Worker Safety and Well-Being through Advanced AI and Emotional Analytics'', Sensors, 24(2), 2024.
- [3] Isik, F., Ozden, G., and Kuntalp, M., ``Importance of Data Preprocessing for Neural Networks Modeling: The Case of Estimating the Compaction Parameters of Soils'', Energy Education Science and Technology Part A: Energy Science and Research, 29(2):871-882, 2012.
- [4] Sharma, S., Pandey, S., and Shah, D., ``Enhancing Medical Diagnosis with AI: A Focus on Respiratory Disease Detection'', Indian Journal of Community Medicine, 48(5), 2023.
- [5] Lakhdari, K., and Saeed, N., ``A New Vision of a Simple 1D Convolutional Neural Networks (1D-CNN) with Leaky-ReLU Function for ECG Abnormalities Classification'', Intelligence-Based Medicine, 6, 2022.
- [6] Akhtar, N., and Ragavendran, U., ``Interpretation of Intelligence in CNN-Pooling Processes: A Methodological Survey'', Neural Comput & Applic, 32:879-898, 2020.
- [7] Mahesh, T. R., Thakur, A., Gupta, M., Sinha, D. K., Mishra, K. K., Venkatesan, V. K., and Guluwadi, S., ``Transformative Breast Cancer Diagnosis Using CNNs with Optimized ReduceLRonPlateau and Early Stopping Enhancements'', International Journal of Computational Intelligence Systems, 17(1), 2024.
- [8] Smith K. Khare, Blanes-Vidal, V., Nadimi, E. S., and Acharya, U. R., ``Emotion Recognition and Artificial Intelligence: A Systematic Review (2014-2023) and Research Recommendations'', Information Fusion, 102, 2024.
- [9] Akinpelu, S., and Viriri, S., ``Deep Learning Framework for Speech Emotion Classification: A Survey of the State-of-the-Art'', IEEE Access, 12:152152-152182, 2024.