

**MAKİNE ÖĞRENMESİ TEKNİKLERİYLE HASTALIK SINIFLANDIRMASI:
RANDOM FOREST, K-NEAREST NEIGHBOUR VE ADABOOST
ALGORİTMALARI UYGULAMASI**

**DISEASE CLASSIFICATION BY MACHINE LEARNING TECHNIQUES: RANDOM
FOREST, K-NEAREST NEIGHBOR AND ADABOOST ALGORITHMS
APPLICATIONS**

1. Doktora Öğrencisi, **Ülkü VERANYURT**

Sağlık Bilimleri Üniversitesi, orcid.org/0000-0003-4838-3373, ulkuveranyurt@gmail.com

2. Yüksek Lisans Öğrencisi, **Ahmet Fatih DEVECİ**

Sağlık Bilimleri Üniversitesi, orcid.org/0000-0002-3044-8397, ahmetfatihdeveci@gmail.com

3. Dr. Öğr. Üyesi, **M. Fevzi Esen**

Sağlık Bilimleri Üniversitesi, orcid.org/0000-0002-3044-8397, fevziesen@gmail.com

4. Yüksek Lisans Öğrencisi, **Ozan VERANYURT**

Bahçeşehir Üniversitesi, orcid.org/0000-0003-3652-2356, ozan.veranyurt@bahcesehir.edu.tr

Makale Gönderim-Kabul Tarihi (22.05.2020-12.08.2020)

Özet

Amaç: Bu çalışmada, sağlık yönetiminde etkinliği sağlamak üzere, hastalıkların doğru olarak teşhisinde makine öğrenmesi tekniklerinin başarısının karşılaştırılması amaçlanmıştır.

Veri Seti ve Yöntem: Çalışmada, Vanderbilt Üniversitesi tarafından çeşitli hastalıkların risk faktörlerinin yaygınlığını anlamak için gerçekleştirilen ve kamuya açık, 390 hastaya ait 15 değişkenden oluşan veri seti kullanılmıştır. Modelin eğitilmesi ve testi amacıyla, veri setinin %70'i eğitim, %30'u test kümelerine bölünmüştür. Random Forest (RF), K-Nearest Neighbour (KNN) ve AdaBoost algoritmaları kullanılarak sınıflandırma performansları karşılaştırılmıştır.

Bulgular: Çalışma sonucunda, RF ve KNN algoritmaları sınıflandırma başarısının %92,30 ve AdaBoost algoritması ile gerçekleştirilen sınıflandırma başarısının ise %90,59 olduğu tespit edilmiştir.

Sonuç: Yapay zekâ ve makine öğrenmesi yöntemlerinin sağlık yönetimi ve hizmetleri alanındaki kullanımı gün geçtikçe artmaktadır. Çalışmamızda, hastalıkların doğru olarak teşhisi amacıyla kullanılan algoritmalarla %90'ın üzerinde doğru sınıflandırma başarısı elde edilmiştir. Bu durum, teşhis ve tedavi süreçlerinde insan kaynaklı

hataları azaltmak ve medikal karar süreçlerine destek amacıyla, makine öğrenmesi tekniklerine başvurulabileceğini göstermektedir.

Anahtar Sözcükler: Hastalık sınıflandırma, sağlıkta makine öğrenmesi, diyabet hastalığı

Abstract

Objective: The aim of this study is to compare the correct classification rates of different machine learning algorithms in the detection of diseases in order to ensure the effectiveness in accurate diagnosis and health management.

Data Set and Method: In our study, we utilized from a research conducted by Vanderbilt University. The aim of the research was to understand the prevalence of risk factors for various diseases. An open access data set of 15 variables of 390 patients belong to that research was used in this study. For the purpose of training and testing of the model, 70% of the data set is divided into training and 30% into test sets. Classification performances were compared using Random Forest (RF), K-Nearest Neighbor (KNN) and AdaBoost algorithms.

Results: As a result of the study, it was observed that the classification success of the RF and KNN algorithms was 92.30% and the classification success of the AdaBoost algorithm was 90.59%.

Conclusion: Artificial intelligence and machine learning methods are used more frequently in the field of health management and services. In our study, the test success was achieved over 90% by using different algorithms. Machine learning techniques can be applied in issues such as reducing human errors in diagnosis and treatment processes and providing support in medical decision making processes.

Keywords: Disease classification, machine learning in health, diabetes disease

GİRİŞ

Sağlık sektöründe tanı ve teşhisin sadece insan gücüyle yönetilmesi mümkün değildir (ADA, American Diabetes Association 2014). Hastalıkların teşhis, tedavi ve rehabilitasyonunun yanı sıra, hastalıkların önlenmesi ve toplumun sağlık düzeylerinin geliştirilmesini de içeren sağlık hizmetlerinin yürütülebilmesi için sağlık yönetimi önem kazanmaktadır. Sağlık hizmetlerindeki değişimin temel sebepleri ise, kronik hastalıkların yaygınlığının artması, hastaların beklentilerinin değişmesi ve giderek artan yaşlı nüfusun evde bakım hizmetleri gibi yeni ihtiyaçların doğması olarak sıralanmaktadır (Hayran, 2012). Bunlara bağlı olarak, yeni teknolojilerin sağlık hizmetleri ve yönetiminde kullanılmaya başlanması da kaçınılmaz görünmektedir.

Sağlık hizmetlerinde süreçler anestezi, görüntüleme, tarama gibi alanlardaki teknolojik gelişmelere bağlı olarak gelişmektedir. Teknolojide inovasyon sağlık hizmetlerini bir dönüşüme sokarken insan faktörü de bu süreçteki gelişmeye açık noktalardan biridir (Thimbleby, 2013). Sağlık yöneticileri, maliyetlerin azalması, verimliliğin artması, bilgi eksikliğinden kaynaklanan personel hatalarını azaltmak veya ortadan kaldırmak ve zamanında ve doğru şekilde hastalık tanımlama yapabilmek için yapay zekâ teknolojileri kullanmaktadır. Yapay zeka teknolojileri insan faktörünün olduğu tüm sağlık hizmetleri ve sağlık yönetimi süreçlerine entegre edilmektedir.

Yapay zekâ teknolojilerinden biri olan makine öğrenmesi teknikleri tıbbi tahminlerde yaygın olarak kullanılmaktadır (Islam vd. 2019; Kononenko, 2001). Makine öğrenmesi teknikleri ile hızlı ve güvenilir şekilde hastalık tahminlemesi yapılabilmektedir. Makine öğrenmesine ait birçok algoritma mevcut olup, problemin kaynağına ve veri sayısına göre hangi algoritmanın kullanılacağına karar



ULUSLARARASI SAĞLIK YÖNETİMİ VE STRATEJİLERİ ARAŞTIRMA DERGİSİ

INTERNATIONAL JOURNAL OF HEALTH MANAGEMENT AND STRATEGIES RESEARCH

Cilt/Volume : 6 Sayı/Issue : 2 Yıl/Year : 2020 ISSN -2149-6161

verilmektedir. Farklı algoritmalar, kullanılan veriye göre özgüllük (specificity) ve duyarlık (sensitivity) değerleri üretebilmektedir (Chubak, Pocobelli ve Weiss, 2012). Kullanılan veriye göre bu değerler göz önünde bulundurularak farklı algoritmalar arasından en iyi uyum sağlayan seçilmektedir. Bu şekilde en iyi algoritma tercih edilerek sonuç iyileştirilmekte olup, zaman maliyeti azaltılmaktadır (Chubak, Pocobelli ve Weiss, 2012).

Çalışmaya konu olan sınıflandırma işlemi için, ölümcül birçok hastalığın oluşumunda önemli rol oynayan ve görülme sıklığı giderek artan ciddi hastalıklardan biri olan diyabet seçilmiştir. Organ kayıplarıyla birlikte, yaşam kalitesini olumsuz yönde etkileyen diyabet, insülini yeterli üretilmediği ya da üretilen insülinin kullanımında meydana gelen bozukluk sonucu kandaki glikoz miktarının artmasıyla ortaya çıkan kronik bir hastalıktır (Glauber ve Karnielli, 2013; ADA, 2010). Amerikan Diyabet Derneği'ne (ADA, 2017) göre, diyabetin tahmini yıllık maliyeti, doğrudan tıbbi maliyetlerde her yıl %5 artışla 245 milyar doları aşmaktadır. Dünya çapında gün geçtikçe artan diyabet prevalansının yanı sıra, her yıl diyabet teşhisi konan hastalar arasında yaklaşık 1,6 milyon ölüm meydana gelmektedir (WHO, 2010). Bununla birlikte, Uluslararası Diyabet Federasyonu (IDF, 2013) tahminlerine göre 2040 yılında dünyada 10 yetişkinden 1'inin diyabet hastası olacağı öngörülmekte olup, diyabet ile ilişkili hastalıkların sağlık harcamalarının 802 milyar ABD dolarını aşacağı belirtilmektedir. Diyabet ve ona bağlı oluşabilen hastalıkların tanı ve tedavi maliyetler yanında, bireyin iş kapasitesinin azalması, ortalama ömrünün azalması ve hasta yakınlarının uğraşlarından dolayı maliyetler de bulunmaktadır (Satman vd. 2013). Bu durum, sağlıkta önleyici faaliyetlerin önemini göstermekte olup, hastalıkların önceden teşhis edilmesinin ne kadar önemli bir maliyet avantajı yaratacağına işaret etmektedir. Maliyetlerin öngörülebilmesi, sağlık yönetimi ve hizmetlerinin sunulmasında avantaj sağlayan bir araç olarak değerlendirildiğinde, hastalıkların doğru olarak teşhis edilmesinde dinamik süreçlerin ve sayısal tekniklerin kullanılması önem kazanmaktadır (Soyiri ve Reidpath, 2013).

Bu çalışmada, diyabet hastalığının erken teşhis edilebilmesi amacıyla, söz konusu hastalığa ilişkin literatürde gösterilen önemli değişkenler seçilerek RF, KNN ve AdaBoost algoritmaları ile hastalığın doğru sınıflandırmasının yapılması amaçlanmıştır. Çalışmada, her bir algoritmanın doğru sınıflandırma performansları karşılaştırılarak en yüksek doğru sınıflandırmanın elde edildiği algoritmalar belirlenmiştir.

LİTERATÜR

Sağlık alanında yapay zekâ ve makine öğrenmesi uygulamalarının kullanımı, tıbbi tanı ve hastalık takip, maliyet tahminleme, görüntüleme analizi, kaynak planlama ve acil durum yönetimi, yapılandırılmamış (unstructured) verinin işlenmesi gibi birçok alt faaliyet alanında gerçekleştirilmektedir (Alonso vd. 2018; Narula vd. 2017; Esteva vd. 2017). Yüksek boyuttaki hasta verilerinin işlevsel hale getirilmesinde de kullanılan yapay zekâ modelleri, veri güvenirliliğinin ve kalitesinin artırılmasında etkin bir rol oynamaktadır (Cichosz, Johansen ve Hejlesen, 2015; Tran vd. 2019). Ancak, yapay zekânın sağlık alanında kullanımına ilişkin; klinik verilerin doğruluğu, verilerin yönetilmesi, verilerin korunmasıyla ilgili yasal ve etik süreçler, sağlık alanında yapay zekâ kullanımını sınırlandırmaktadır (Char, Shah ve Magnus, 2018). Veri gizliliği ve korunması birçok uygulamada bir sorun teşkil ederken, konuya ilişkin etik süreçler her ülkede farklılık göstermektedir (Celebi ve Inal, 2019).

Diyabet hastalığının tespiti konusunda yapay zekâ uygulamalarıyla gerçekleştirilen farklı çalışmalar mevcuttur. Mujumdar ve Vaidehi (2019) makine öğrenmesiyle diyabet hastalığının tespitinde doğruluk oranını arttırmak için sağlık yönetim sistemindeki büyük veri kaynakları kullanmıştır.

Çalışmada, mevcut hasta veri kümesi bu büyük veri kaynağından alınan yeni parametreler ile zenginleştirilmiş, bu sayede hastalık tanısının konmasında makine öğrenmesinin doğruluk oranı arttırılmıştır. Yapılan farklı bir çalışmada ise, Hindistan örneklemindeki hastalara ilişkin bir veri kümesi kullanılarak gözetimli (supervised) makine öğrenmesi algoritmaları ile hastalık tespiti üzerinde çalışılmıştır (Kaur ve Kumari, 2018). Sağlık yönetimi süreçlerinde diyabet hastalığının yönetilmesi ile ilgili yapılan bir başka çalışmada ise, genel bir inceleme yapılmış olup, hastalığın tahminlenmesi ve sonrasında oluşan komplikasyonlar, hastaların genetik geçmişleri, tedavi ve yönetim süreçleri ile ilgili yapılmış olan çalışmalar incelenmiştir. Çalışmada ayrıca, makine öğrenimine dayalı çalışmaların neredeyse %85'inin sınıflandırma ve tahmin için denetimli öğrenme algoritmaları kullandığı vurgulanmaktadır. Buna göre, en çok tercih edilen algoritma Destek Vektör Makinesi (Support Vector Machine) olmuştur.

Makine öğrenmesi teknikleri sağlık bilimleri alanında en yaygın olarak tahminleme, tanı, hastalık sonrası komplikasyonların belirlenmesinde kullanılmakta olup, zamandan ve iş yükünden tasarruf yapılarak hastalara daha kaliteli sağlık hizmeti verilmesi amaçlanmaktadır (Kavakiotis vd. 2017). Woldaregay ve arkadaşları (2019), diyabetli hastalar için, sağlık yönetimi süreçlerinde hastanın konforunu sağlama ve tedavi süreçlerini iyileştirmeye yönelik hasta yaklaşımı ve yapay zekâ tabanlı hasta takip sistemi önermiştir.

Makine öğrenmesi teknikleri çözümlerine ilişkin sağlık yönetimi süreçlerine bir başka örnek ise, sağlık sigortası sağlayıcıların ön onay süreçlerine makine öğrenme entegre edilerek iyileştirmesi ve hızlandırılmasıdır (Araújo vd. 2016). Sağlık sigortası sağlayıcılarının ön onay süreçlerinde yanlış girilen form, kapsam dışı tedavi veya hizmet gibi taleplerin ön incelemede hızlı ve doğru değerlendirilebilmesi amacıyla, makine öğrenmesi teknikleriyle otomatizasyonu mümkün kılınmaktadır. Sağlık hizmetlerinde makine öğrenmesinin yarattığı fırsatlarla ilgili bir başka çalışmada ise, hastane verileri kullanılarak hastaları risk gruplarına göre sınıflandıran bir model önerilmiş, bu şekilde hasta takibi ve ilgili aksiyonlar konusunda bir iyileşme sağlanmıştır (Parikh, Kakad ve Batesi 2016). Bir diğer çalışmada ise, gereken tedavi ve adımları belirlemek için yüksek maliyetli ve riskli hastaları belirleyen bir yapay zekâ modeli önerilmiştir (Bates vd. 2014). Mercado ve arkadaşlarının (2017) yaptığı bir çalışmada diyabet tanısı ve sınıflandırılması için Hoefding Tree algoritmasını kullanmıştır. Sağlık yönetiminin bir parçası olan hastane yönetimi ile ilgili farklı bir bakış açısı Rodriguez ve arkadaşları (2019) tarafından sunulmuştur. Söz konusu çalışmada, diyabet ve benzeri kronik hastalık sebebiyle yatış yapan hastaların maliyetleri üzerinden bir makine öğrenmesi modeli oluşturulmuş olup, model ile hastalık maliyetlerinin tahminlenmesi amaçlanmıştır.

Bu çalışmada kullanılacak olan, RF, KNN ve Adaboost algoritmalarının etkinliğine ilişkin literatürde farklı alanlarda çalışmalar bulunmaktadır. Söz konusu çalışmalarda, algoritmaların sınıflandırma performansları karşılaştırılmamış olup, her bir algoritmanın güçlü ve zayıf yönleri farklı çalışmalarda tartışılmıştır. Vishwanath ve arkadaşları (2020), farelerde beyin zedelenmesi sınıflandırmasını EEG (electroencephalogram) verileri kullanarak, RF, KNN ve Evrimsel Sinir Ağı (CNN) algoritmaları ile gerçekleştirmiş olup söz konusu çalışmada en yüksek doğru sınıflandırma skoru CNN ile elde edilmiştir. Kalra ve arkadaşları (2020), CT ve MRI kullanılan protokollerin kalite ve etkinliğini doğal dil işleme (NLP) tabanlı bir makine öğrenmesi sınıflandırıcı ile arttırmayı denemiştir. Buna göre, KNN ve RF makine öğrenmesi sınıflandırıcıları kullanılmış olup, algoritmaların birbirlerine karşı herhangi bir üstünlüğü tespit edilememiştir. Wang ve arkadaşları (2020), akciğer kanserinin alt tiplerini belirlemede KNN algoritmasını kullanmıştır. Çalışmada, solunum fonksiyon testi sonuçları kullanılarak hastalardaki akciğer kanseri tipleri sınıflandırılmış olup, KNN algoritmasının yüksek oranda doğru sınıflandırma başarısı gösterdiği sonucuna ulaşılmıştır. Liu ve arkadaşları (2020) ise aort

diseksiyonu (AD) taraması için farklı makine öğrenmesi algoritmaları ile bir metod geliştirmeye çalışmıştır. Söz konusu çalışmada AdaBoost, SmoteBagging, EasyEnsemble ve CalibratedAdaMEC algoritmaları kullanılmıştır. Çalışma sonucunda algoritmaların etkinliği tespit edilmiş olup, söz konusu algoritmalarla gerçekleştirilen sınıflandırmada yanlış teşhis oranlarının %25'in altında olduğu belirlenmiştir.

VERİ SETİ VE METOD

Bu çalışmada, Vanderbilt Üniversitesi Biyoistatistik Departmanı tarafından açık erişim olarak sağlanan diyabet hastası 390 kişiye ait 15 değişkenden oluşan veri seti kullanılmıştır. Çalışmada kullanılan veri kamuya açık olup, 1993-1994 yılları arasında toplanmıştır. Bu nedenle bir etik kurul izni gerektirmemektedir. Veri seti, Random Forest, K-Nearest Neighbour ve AdaBoost algoritmaları kullanılarak sınıflandırılmıştır.

Makine öğrenmesi, iterasyonlar ve veri üzerindeki ortak desenleri herhangi bir varsayım olmadan öğrenmeye çalışmaktadır. Makine öğrenmesi çeşitleri Tablo 1'de özetlenmiştir. Bu çalışmada gözetimli olarak RF ve AdaBoost, gözetimsiz olarak ise KNN makine öğrenmesi tekniği seçilerek, probleme ilişkin doğru sınıflandırmada uygun algoritmanın bulunması amaçlanmıştır.

RF algoritması, karar ağaçlarının birleşiminden oluşan bir algoritma olup, kullanılan karar ağaçları arasında doğruluğu ve bağımsızlığı en yüksek ağaçlar tercih edilmektedir. Ağaçların her birisi veri setindeki özelliklere göre dallanmaktadır. Dallanma bu özelliklerde karar noktalarına bağlıdır (Breiman, 2001). Karar ağacı kullanan modellerdeki en büyük problem, verinin az olması durumunda aşırı oranda uyumun (over fitting) oluşmasıdır (Liao, Ju ve Zou, 2016). Random Forest sınıflandırıcısı, veri kümesinin çeşitli alt örneklerine bir dizi Karar Ağacı Sınıflandırmasına (Decision Tree Classifier) uyan ve tahmin doğruluğunu iyileştirmek ve aşırı uyumu kontrol etmek için ortalamayı kullanan bir meta tahmincidir. Alt örnek boyutu her zaman orijinal giriş örneklem boyutuyla aynıdır.

KNN algoritması ise gözetimsiz bir algoritma olup, veri içerisindeki noktaların kendisine en yakın noktalar ile karşılaştırılarak veri kümeleri oluşturulması hedeflenmektedir. Komşu sayısı parametresiyle, karşılaştırmanın veriye en yakın kaç komşu ile yapılacağı belirlenmektedir (Jiang ve Zhou, 2016). Algoritmayı bilgisayar ortamında çalıştırabilmek için tüm veri setine ihtiyaç duyulmakta olup, tüm veri setinin ana bellekte öğrenme süreci boyunca tutulması gerekmektedir. Veriler arasındaki mesafeyi ölçmek için Öklit, Manhattan gibi mesafe ölçütleri kullanılmaktadır (Cover ve Hart, 1967). KNN 'de amaç, yeni noktaya en yakın önceden belirlenmiş sayıda eğitim örneği bulmak ve etiketi bunlardan tahmin etmektir. Örnek sayısı kullanıcı tanımlı bir sabit (k-en yakın komşu öğrenimi) olabilir veya noktaların yerel yoğunluğuna (yarıçap temelli komşu öğrenimi) bağlı olarak değişebilir. Basitliğine rağmen, en yakın komşular el yazısı rakamlar ve uydu görüntüsü sahneleri de dâhil olmak üzere çok sayıda sınıflandırma ve regresyon probleminde başarılı olarak kullanılmaktadır. Parametrik olmayan bir yöntem olduğundan dolayı, karar sınırının çok düzensiz olduğu sınıflandırma durumlarında genellikle başarılı olmaktadır.

Tablo 1. Makine Öğrenmesi Türleri

Metod	Tanım
Gözetimli	Yaygın olarak kullanılan öğrenme biçimidir. Veri sırasıyla eğitim ve test kümelerine ayrılır. Yapılacak işleme göre veri önceden işaretlenmelidir. Regresyon, tahminleme, sınıflandırma gibi işlemlerde kullanılır. Yapay sinir ağları, Random Forest gibi algoritmalar örnek gösterilebilir.
Gözetimsiz	Bu öğrenme biçiminde modele eğitim için herhangi bir sınıf ya da sayısal bilgi verilmez. Model veri içindeki ortak noktalar üzerinden sonuç üretmeye çalışır. K-Nearest Neighbour (KNN), hierarchical clustering örnek verilebilir.
Yarı gözetimli	Veri içinde işaretlenmiş sonuç bilgisi olan ya da olmayan bölümlere ayrılır. Tam olarak sınıflandırılmamış bir veri kullanılır. Ses algılama buna örnek verilebilir.
Pekiştirmeli	Davranışsal psikolojideki ödül prensibine dayanır. Karar verme mekanizması alınan sonuca göre en yüksek ödülü veren seçeneği öğrenir. Günümüzde robotik, yapay zekâda kişiselleştirme, medikal görüntü işleme gibi alanlarda kullanılmaktadır.

Kaynak: (Shameer vd. 2017)

Adaboost algoritması, Freund ve Schapire (1997) tarafından önerilmiş olup, çeşitli zayıf algoritmaların bir araya getirilerek karma ve daha güçlü bir eğitim modeli sunan bir makine öğrenmesinin gerçekleştirilmesi amaçlanmaktadır. İçerisinde kendisinden daha önce geliştirilmiş olan farklı makine öğrenmesi algoritmalarını birleştirerek ya da karşılaştırarak öğrenme konusunda en iyi sonucun elde edilmesi hedeflenmektedir (Freud ve Schapire, 1997).

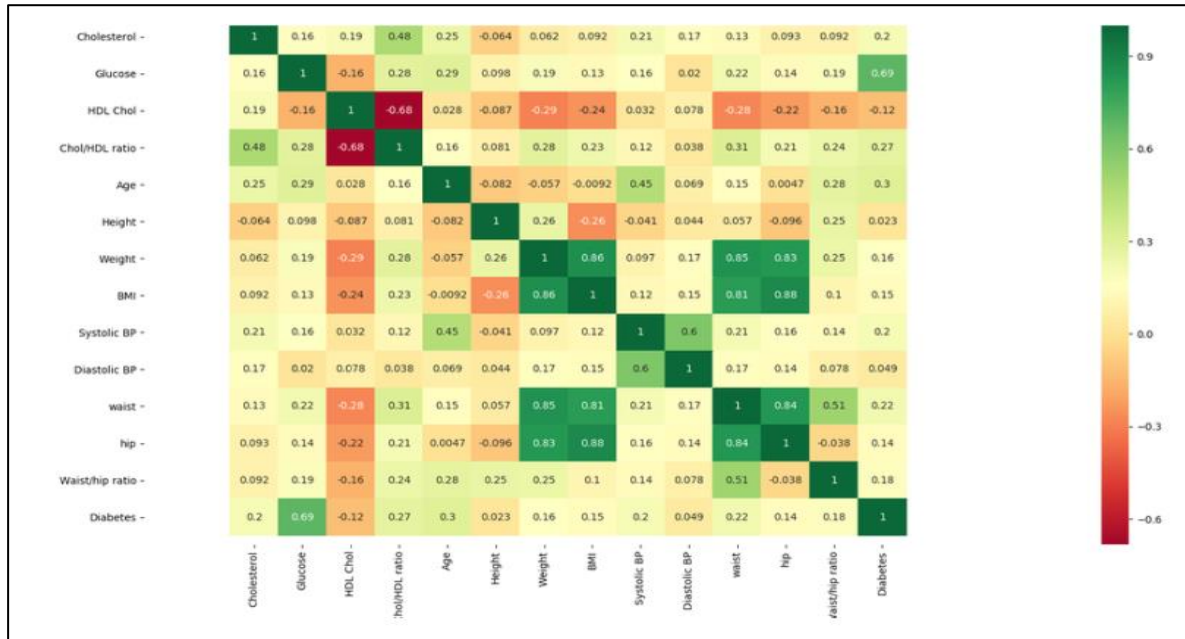
Bu çalışmada kullanılan veriler obezite, diyabet ve kardiyovasküler risk faktörlerinin yaygınlığının tespitine ilişkin farklı hastalık gruplarında bulunan kişilerden elde edilmiştir. Veri kümesinde, yukarıda belirtilen hastalıklarla ilişkili 15 farklı değişken bulunmaktadır. Veri setinin %70'i eğitim, %30'u ise test verisi olarak ayrılmıştır. Veri setine kişilerin yaş, boy, ağırlık, vücut kitle endeksi (BMI), sistolik kan basıncı, diastolik kan basıncı, bel, bel/kalça oranı, kalça, kolesterol, HDL kolesterol, kolesterol/HDL oranı ve glikoz verileri dâhil edilmiştir (Tablo 2). ADA'ya göre yaş, cinsiyet, boy uzunluğu, kilo, Vücut Kitle İndeksi (BMI), sistolik kan basıncı, diastolik kan basıncı, bel, kalça, bel/kalça oranı, kolesterol, HDL kolesterol ve kolesterol/HDL kolesterol oranı gibi parametreler diyabet hastalarında; kardiyovasküler hastalıklar için risk faktörleridir (ADA, 2019). Ayrıca obeziteyi gösteren yaş, cinsiyet, boy, ağırlık, BMI, bel, kalça, bel/kalça oranı gibi değişkenler diyabet hastalığı için önemli belirteçlerdir (ADA, 2019).

Veri setine ilişkin temel istatistikî ölçütler Tablo 2'de gösterilmiştir. Örneklemin yaş aralığı 19-92 yıl gibi geniş bir grubu kapsamakta olup, yaş ortalaması 46,77 olarak bulunmuştur. Çalışmada değerlendirilen hastaların %41.53'ü erkek (n=162), %58.46'sı kadındır (n=228). Ayrıca, hastaların %15.38'inin (n=60) diyabet tanısı bulunmaktadır.

Tablo 2: Temel İstatistiksel Ölçütler

Değişken	Min-Max	Ortalama ± Std. Sapma
Yaş	19-92	46,77 ± 16,43
Boy	132-193 cm	167,50 ± 9,93
Ağırlık	45-147 kg	80,46 ± 18,32
Vücut Kitle Endeksi (BMI)	15,2 – 55,8 lbs/in ²	28,77 ± 6,60
Sistolik Kan Basıncı	90-250 mmHg	137,13 ± 22,85
Diastolik kan basıncı	48-124 mmHg	83,28 ± 13,49
Bel	66-142 cm	96,16 ± 14,63
Bel/Kalça oranı	0,68-1,14	0,88 ± 0,07
Kalça	76-162 cm	109 ± 14,37
Kolesterol	78-443 mg/dl	207,23 ± 44,66
HDL Kolesterol	12-120 mg/dl	50,26 ± 17,27
Kolesterol/HDL oranı	1,5-19,3	4,52 ± 1,73
Glukoz	48-385 mg/dl	107,33 ± 53,79

Değişkenler arasındaki korelasyon Şekil 1’de gösterilmiştir. Renk skalasına göre yeşile yakın olan sütunlar arasında korelasyon ilişkisi yüksek iken, kırmızıya doğru olan renklerde ise korelasyon ilişkisi giderek düşmektedir.

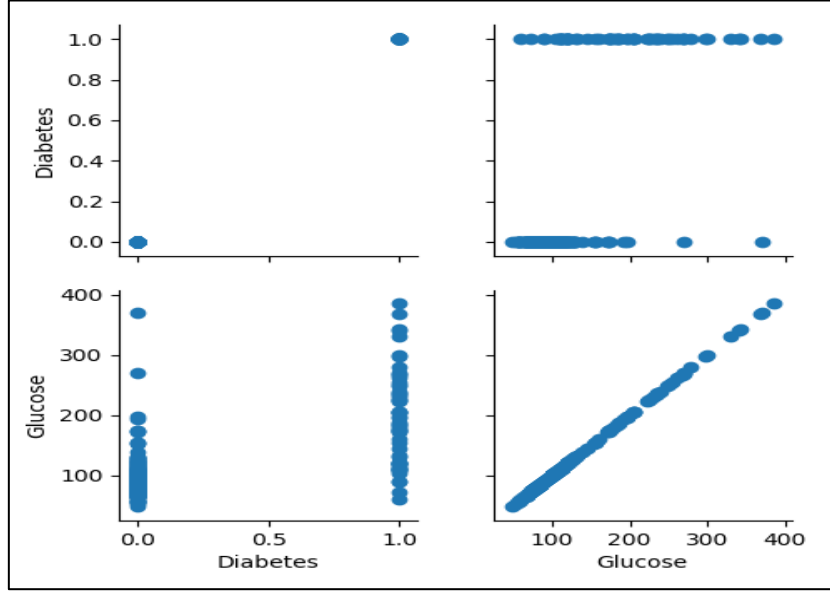


Şekil 1. Korelasyon Matrisi

Veri setinde hastanın diyabet durumu yani bağımlı değişken ile diğer sütunlar arasında en yüksek korelasyona glikoz değişkeninin sahip olduğu görülmektedir. Ayrıca, bağımlı değişken (diyabet

hastası / diyabet hastası değil) ile diğer değişkenler arasında 0,5'in üzerinde bir korelasyon tespit edilememiştir. Bu durum, hiç bir değişkenin tahminlemede bağımsız olarak kullanılamayacağını göstermektedir.

Hastaların glikoz değerleri ve diyabet hastalığına sahip olma durumu arasındaki ilişki ise Şekil 2'de gösterilmiştir. Buna göre, veri tam bir bölümlenme göstermekte olup, her iki değişken arasında korelasyon bulunmaktadır.

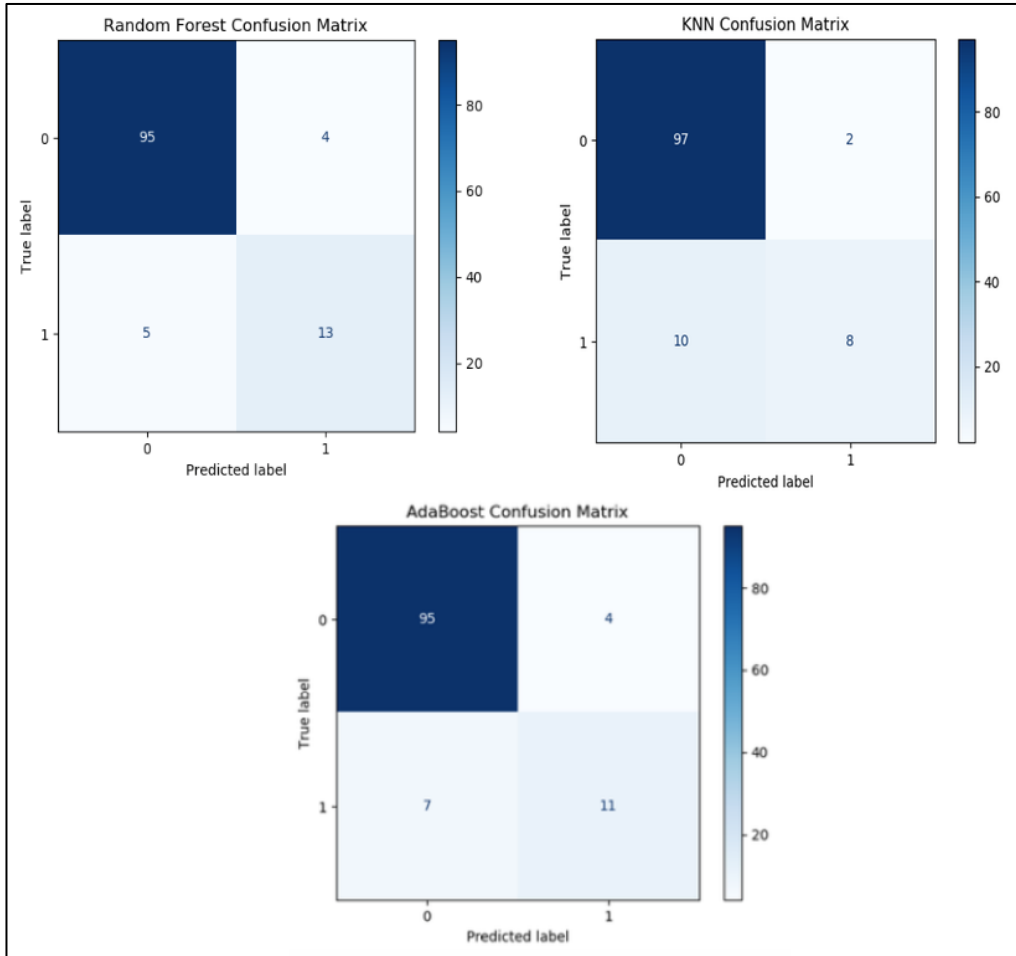


Şekil 2. Glikoz değeri – Diyabet hastalığı ilişkisi

RF, KNN ve AdaBoost algoritmaları kullanılarak veri kümesi eğitim ve test kümelerine bölünerek her bir algoritma için öğrenme gerçekleştirilmiştir. Her bir algoritma için kümelendirme ve doğrulama işlemi 10 defa tekrarlanmış olup, ortalama değerler göz önünde bulundurulmuştur. Söz konusu işlemler veri üzerinde ölçeklendirme yapılmaksızın uygulanmıştır.

BULGULAR

Normalizasyon yapılmadan gerçekleştirilen sınıflandırmaya göre, RF ve KNN algoritmalarıyla yapılan doğru sınıflandırma oranının %92,30 olduğu ve birbirine benzer sonuçlar ürettiği, AdaBoost algoritması ile yapılan doğru sınıflandırma başarısının ise %90,59 olduğu gözlemlenmiştir. Algoritmaların ölçeklendirilmemiş veri üzerinde tek seferlik çalışmaları değerlendirildiğinde, RF algoritmasının Şekil 3'de gösterildiği gibi 117 test verisinden 108'ini doğru olarak tahminlediği görülmektedir. Söz konusu test için hiperparametre olarak ağaç sayısı, n=100 olarak belirlenmiş olup, daha düşük değerlerde veri doğruluğunda azalma gözlemlenmiştir. Aynı şekilde, KNN algoritması sonuçlarına göre 117 test verisinde 105 kişi için doğru tahminleme yapılmıştır. KNN algoritması için 15 komşu sayısına kadar model eğitilip 3 komşu sayısında en yüksek doğruluk elde edilmiştir. AdaBoost algoritması sonucunda ise, 117 test verisinde 106 kişi için doğru tahminleme yapılmış olup, algoritma için maksimum tahminci sayısı 10 olarak seçilmiştir. Öğrenme sıklığı 0,1 olarak belirlenmiş olup, en yüksek doğruluğa ulaşıldığında model eğitimi durdurulmuştur.



Şekil 3. Confusion Matrisleri

Çalışma sonucunda, diyabetin doğru sınıflandırmasında RF ve KNN algoritmaları %92,30'lik başarı sağlamıştır. Yanlış teşhis oranı açısından en düşük hata oranı yine RF ve KNN algoritmaları ile sağlanmıştır.

TARTIŞMA

Çalışmamızda hastalık sınıflandırma başarısı KNN ve Random Forest algoritmaları ile daha yüksek bulunmuş olup, söz konusu algoritmalarla yüksek doğrulukla sınıflandırma yapılabildiği anlaşılmaktadır. Çalışma sonuçlarının, Maniruzzaman ve arkadaşları (2018) ve Chen ve Pan (2018) tarafından yapılan hastalık tahminlenmesine yönelik çalışmalar ile uyumlu olduğu görülmektedir.

Yapay zekâ otomotiv, pazarlama, hizmet vb. birçok sektörde kullanılmakla birlikte sağlık sektöründeki önemi ve kullanımı her geçen gün artmaktadır. Yapay zekâ teknolojileri kullanılarak sağlık hizmetleri ve yönetiminde erken tanı, tedavi ve operasyonel süreçlerin planlanmasına ilişkin birçok konuda uygulama alanı bulunmaktadır. Özellikle ömür boyu tedavi gerektiren ve maliyeti yüksek, kronik hastalıkların erken teşhisinde makine öğrenmesi algoritmalarının kullanım alanının

genişletilmesi, sağlık hizmetlerinin sunumunda kalite ve maliyet açısından optimizasyonu mümkün kılacaktır. Bu doğrultuda, makine öğrenmesi tekniklerinin mobil yada web üzerinden uzman bir sistemle entegre bir şekilde geliştirilmesi; hastalıkların erken safhada tespit edilerek, söz konusu hastalıkların sebep olduğu sekonder hastalıkların (körlük, organ amputasyonu, böbrek hastalığı vb.) oluşmadan engellenmesi sağlanabilir. Bu yolla, hem bireylerin hayat kalitesi artırılırken, hem de hastalığın sağlık yönetimi açısından maliyet ve iş yükünün azaltılması mümkün olabilir. Oluşturulan model ve kullanılan analizlerin güvenilirliği açısından, çoklu validasyon veya farklı veri kümeleri kullanılarak ve veri hacmi, boyutu artırılarak hastalıkların teşhis ve tedavi hızını da içerisine alan sağlık yönetimi süreci optimizasyonu çalışmalarına öncelik verilebilir.

KAYNAKÇA

ADA. (2010). Diagnosis and classification of es mellitus. *es Care*, 33(1): 62-69.

ADA. (2014). Standards of medical care in. *es Care*, 37(1): 14-80.

ADA (2017). Economic Costs of Diabetes in the U.S. in 2017. ADA. doi: <https://doi.org/10.2337/dci18-0007>.

ADA. (2019). Cardiovascular Disease and Risk Management: Standards of Medical Care in Diabetes. *Diabetes Care*, 42 (1):103-123 | <https://doi.org/10.2337/dc19S010> (02.04.2020).

Alonso, DH., Wernick, MN., Yang, Y., Germano, G., Berman, DS., Slmoka, P. (2018). Prediction of cardiac death after adenosine myocardial perfusion SPECT based on machine learning. *J Nucl Cardiol*. <https://doi.org/10.1007/s12350-017-0924-x> (02.02.2020).

Araújo F.H.D. et al. (2016). Using machine learning to support healthcare professionals in making pre authorization decisions. *International Journal of Medical Informatics*, 94:1-7.

Bates, DW., Saria, S., Ohno-Machado, L., Shah, A., Escobar, G., (2014). Big data in healthcare: using analytics to identify and manage high-risk and high-cost patients. *Health Aff*, 33: 1123-1131.

Breiman, L. (2001). Random forest. *Mach. Learn*, 45: 5-32. doi: 10.1023/A:1010933404324.

Char, DS., Shah, NH., Magnus, D. (2018). Implementing Machine Learning in Health Care Addressing Ethical Challenges. *N. Engl. J. Med.*, 378: 981-983.

Chen, P. and Pan, C. (2018). Diabetes classification model based on boosting algorithms. *BMC*, 19:109 <https://doi.org/10.1186/s12859-018-2090-9> (14.03.2020).

Celebi, V., Inal, A. (2019). Problem of Ethics in the Context of Artificial Intelligence. *The Journal of International Social Research*, 12, 66.

Chubak, J., Pocobelli, G., Weiss, NS. (2012). Trade-offs between accuracy measures for electronic healthcare data algorithms. *J Clin Epidemiol*, 65(3):343-349.e2.

Cichosz, SL., Johansen, MD., Hejlesen, O. (2015). Toward big data analytics: review of predictive models in management of es and its complications. *J es Sci Technol*, 10(1):27-34.



ULUSLARARASI SAĞLIK YÖNETİMİ VE STRATEJİLERİ ARAŞTIRMA DERGİSİ

INTERNATIONAL JOURNAL OF HEALTH MANAGEMENT AND STRATEGIES RESEARCH

Cilt/Volume : 6 Sayı/Issue : 2 Yıl/Year : 2020 ISSN -2149-6161

Cover, T. and Hart, P. (1967). Nearest neighbor pattern classification,"Information Theory, IEEE Transactions, 13: 21-27.

Esteva, A., Kupre, B., Novoa, RA., et al. (2017). Dermatologist-level classification of skin cancer with deep neural networks. Nature; 542:115–8

Freund, Y and Schapire, RE. (1997). A decision-theoretic generalization of on-line learning and an application to boosting. Journal of Computer and System Sciences, 55(1):119–139.

Glauber, H., Karnieli, E. (2013) Preventing type 2 es mellitus: a call for personalized intervention. Perm J, 17(3): 74-9

Hayran, O. (2012). Sağlık Yönetimi Yazıları. Sage Yayıncılık: Ankara.

IDF. Atlas. (2013). 6th edition, <http://www.idf.org/esatlas> (14.03.2020).

Islam, T., Raihan, M., Farzana, F. et al. (2019.) An Empirical Study on es Mellitus Prediction for Typical and Non-Typical Cases using Machine Learning Approaches. 10th ICCCNT 2019. Kanpur, India.

Jiang, Y and Zhou, ZH. (2004). Editing training data for kNN classifiers with neural network ensemble. Lect. Notes Comput. Sci. 3173: 356–361. doi: 10.1007/978-3-540-28647-9_60.

Kalra, A., Chakraborty, A., Fine, B., Reicher, J. (2020). Machine Learning for Automation of Radiology Protocols for Quality and Efficiency Improvement. J Am Coll Radiol. doi: 10.1016/j.jacr.2020.03.012.

Kaur, H., Kumari, V. (2018). Predictive modelling and analytics for diabetes using a machine learning approach. Applied Computing and Informatics <https://doi.org/10.1016/j.aci.2018.12.004> (05.03.2020).

Kavakiotis, I. et al. (2017). Machine Learning and Data Mining Methods in Diabetes Research. Computational and Structural Biotechnology, 15: 104–116.

Kononenko, I. (2001). Machine learning for medical diagnosis: history, state of the art and perspective. Artificial Intelligence in Medicine, vol. 23, no. 1, pp. 89-109

Liao, Z., Ju, Y., and Zou, Q. (2016). Prediction of G protein-coupled receptors with SVM-Prot features and random forest. Scientifica, 8309253. doi: 10.1155/2016/8309253.

Liu, L., Zhang, C., Zhang, G. et al. (2020). A study of aortic dissection screening method based on multiple machine learning models. J Thorac Dis, 12(3):605-614. doi: 10.21037/jtd.2019.12.119.

Maniruzzaman, M., Rahman, MJ., Al-Mehedi Hasan, M. et al. (2018). Accurate es Risk Stratification Using Machine Learning: Role of Missing Value and Outliers. J Med Syst, 10;42(5):92. doi: 10.1007/s10916-018-0940-7.



ULUSLARARASI SAĞLIK YÖNETİMİ VE STRATEJİLERİ ARAŞTIRMA DERGİSİ

INTERNATIONAL JOURNAL OF HEALTH MANAGEMENT AND STRATEGIES RESEARCH

Cilt/Volume : 6 Sayı/Issue : 2 Yıl/Year : 2020 ISSN -2149-6161

- Mercaldo, F., Nardone, V., Santone, A. (2017). Diabetes Mellitus Affected Patients Classification and Diagnosis through Machine Learning Techniques. *Procedia Computer Science*, 112: 2519-228.
- Mujumdar, A., Vaidehi, V. (2019). Diabetes Prediction Using Machine Learning Algorithms. *Procedia Computer Science*, 165: 292–299.
- Narula, S., Shameer, K., Salem Omar, AM., Dudley, JT., Sengupta, PP. (2017) Reply: Deep learning with unsupervised feature in echocardiographic imaging. *J Am Coll Cardiol*;69:2101–2.
- Parikh, R.B., Kakad, M., Bates, DW. (2016). Integrating predictive analytics into high-value care: the dawn of precision delivery. *JAMA*, 315: 651-652.
- Rodriguez, G. et al. (2019). Predicting Healthcare Costs of Diabetes Using Machine Learning Models. Elsevier Inc., doi: <https://doi.org/10.1016/j.jval.2019.09.903> (05.04.2020).
- Satman, I., Omer, B., Tutuncu, Y., Kalaca, S., Gedik, S., Dinssccag N, Karsidag, K. & TURDEP-II Study Group. (2013). Twelve-year trends in the prevalence and risk factors of es and prees in Turkish adults. *Eur J Epidemiol*, 28(2):169-180.
- Shameer K, Johnson KW, Glicksberg BS, Dudley JT, Sengupta PP. (2017). Machine learning in cardiovascular medicine: Are we there yet? <https://doi.org/10.1136/heartjnl-2017-311198> (05.04.2020).
- Soyiri, NI., Reidpath, DD. (2013). An overview of health forecasting. *Environ Health Prev Med* 18:1–9. DOI 10.1007/s12199-012-0294-6.
- Thimbleby H. (2013). Technology and the future of healthcare. *Journal of Public Health Research*; 2:e28.
- Tran, BX., Latkin, CA., Giang, VT., et al. (2019). The Current Research Landscape of the Application of Artificial Intelligence in Managing Cerebrovascular and Heart Diseases: A Bibliometric and Content Analysis. *Int. J. Environ. Res. Public Health*, 16:2699.
- Vishwanath, M., Jafarlou, S., Shin, I. et al. (2020). Investigation of Machine Learning Approaches for Traumatic Brain Injury Classification via EEG Assessment in Mice. *Sensors (Basel)*, 20(7). doi: 10.3390/s20072027.
- Wang, C., Long, Y., Li, W. et al. (2020). Exploratory study on classification of lung cancer subtypes through a combined K-nearest neighbor classifier in breathomics. *Sci Rep*, 3;10(1):5880. doi: 10.1038/s41598-020-62803-4.
- WHO. (2020). Diabetes. <https://www.who.int/health-topics/diabetes> (14.03.2020).
- Woldaregaya, AZ. et al. (2019). Data-driven modeling and prediction of blood glucose dynamics: Machine learning applications in type 1 diabetes *Artificial Intelligence in Medicine*, 98: 109–134.