





# Düzce Üniversitesi Bilim ve Teknoloji Dergisi

Araştırma Makalesi

## Parkinson Hastalığı Teşhisi İçin Makine Öğrenmesi Tabanlı Yeni Bir Yöntem

 Sadullah ESMER <sup>a\*</sup>,  Muhammed Kürşad UÇAR <sup>b</sup>,  İbrahim ÇİL <sup>c</sup>,  
 Mehmet Recep BOZKURT <sup>b</sup>

<sup>a</sup> Elektrik Elektronik Mühendisliği Bölümü, Mühendislik Fakültesi, Bolu Abant İzzet Baysal Üniversitesi, Bolu, TÜRKİYE

<sup>b</sup> Elektrik Elektronik Mühendisliği Bölümü, Mühendislik Fakültesi, Sakarya Üniversitesi, Sakarya, TÜRKİYE

<sup>c</sup> Endüstri Mühendisliği Bölümü, Mühendislik Fakültesi, Sakarya Üniversitesi, Sakarya, TÜRKİYE

\* Sorumlu yazarın e-posta adresi: sadullahesm@hotmail.com

DOI: 10.29130/dubited.688223

### ÖZET

Parkinson hastalığı (PH), dopamin üreten beyin hücrelerinin ölmesiyle ya da zarar görmesiyle ortaya çıkan bir beyin hastalığıdır. Böyle bir durumda, beyin normal fonksiyonlarını yerine getiremez. PH, konuşma, yürüme ve yazma gibi insan hareketlerini olumsuz olarak etkiler. Bu hastalığın teşhisinde detaylı tıbbi öykü, geçmiş tedavi öyküsü, fiziksel testler ve bazı kan testleri ile beyin filmleri istenilmektedir. Bu işlemler maliyetli ve meşakkatli olabildiği için daha az maliyetli ve daha kolay yapılabilen teşhis bu noktada önem kazanmaktadır. Bu çalışmada doktorun kararına destek olabilmesi için 252 bireyden alınan ses verileri ile PH'nin teşhis edilebilmesi amaçlanmıştır. Verilerden daha iyi sonuç alabilmek için bazı ön işlemler uygulanmıştır. Verilerde dengeleme işlemi yapılmış ve sistematik örnekleme metodu ile işleme alınacak veriler belirlenmiştir. Öznitelik seçme algoritması ile özniteliklerin etiket üzerindeki etki gücü hesaplanıp bazı veri grupları oluşturulmuştur. Sınıflandırma algoritmalarından Karar ağacı, Destek Vektör Makineleri ve K En Yakın Komşu Algoritması kullanılıp, performans değerlendirme kriterleri - bunlar; Doğruluk Oranı, Duyarlılık, Özgünlük, F-Ölçümü, Kappa, Auc - değerlendirilmiştir. En yüksek performans değerine sahip veri grubu ve kullanılan sınıflandırma algoritması belirlenip model oluşturulmuştur. Model en ilgiliden en ilgisize doğru sıralanmış veri setinin %45'inden ve Destek vektör makineleri algoritması kullanılarak oluşturulmuştur. Performans kriterlerinde %85 doğruluk oranı ve diğer kriterlerde de olumlu sonuçlar elde edilmiştir. Böylece PH olma ihtimali olan bireyin ses kayıtlarından oluşturulan veri seti ve uygulanan model yardımı ile doktora tıbbi karar destek sağlanacağı anlaşılmıştır.

**Anahtar Kelimeler:** Tıbbi Karar Destek Sistemleri, Sınıflandırma Algoritmaları, Parkinson Hastalığı, Makine Öğrenmesi

## A New Method Based on Machine Learning for the Diagnosis of Parkinson's Disease

### ABSTRACT

Parkinson's disease (PD) is a brain disease caused by death or damage of dopamine producer brain cells. In such case, the brain can not perform its normal functions. PD negatively affects human movements such as speech, walking and writing. In the diagnosis of this disease, detailed medical history, history of treatment, physical tests

and some blood tests and brain films are required. Because these operations can be costly and difficult, less costly and easier making of the diagnosis has such important in this subject. In this study, it was aimed to diagnose PD with voice data from 252 individuals to support the doctor's decision. In order to get better results, some pre-treatments were applied. The datas were balanced and datas that taken to treatments with systematic sampling method were determined. With the feature selection algorithm, some data groups were created by calculating the effect of the attributes on the label. Of the classification algorithms; Decision tree, Support Vector Machines and K Nearest Neighbor Algorithm are used and Performance evaluation criteria such as Accuracy, Sensitivity, Specificity, F-Measurement, Kappa, Auc - were evaluated. The data group with the highest performance value and the used classification algorithm were determined and model was created by using support vector machines algorithm and from 45% of data set that was sorted from the most effective to the most ineffective. The performance criterias has an accuracy of 85% besides in other criterias positive results were earned. Thus, it was understood that medical decision support to the doctor would be provided with the help of applied model and data set formed by sound recordings of the individuals who possibly been PD

*Keywords: Clinical Decision Support Systems, Classification Algorithms, Parkinson's Disease, Machine Learning*

## I. GİRİŞ

Parkinson hastalığı (PH) yaygın bir nörolojik rahatsızlıktır. Hastalığın tipik semptomları ilk olarak 19. yüzyılda İngiliz bilim adamı James Parkinson (1755-1824) tarafından detaylı bir şekilde açıklanmıştır. Hastalık adını James Parkinson'dan almaktadır. Bu hastalığın görülme sıklığı 1000'de 1'dir. PH, 60 yaşın üzerindeki bireylerin %1'inde görülürken, 85 ve üzeri yaşlarda bu oran %5'lere çıkmaktadır [1]. PH, harekette yavaşlık (bradykinesia), titreme ve kasılma olarak karakterize edilmektedir [2]. Bunlara ek olarak uyku bozukluğu, depresyon belirtileri ve konuşma bozukluğu görülmektedir [3]. Konuşma bozukluğu kısık sesle konuşma, donuk konuşma, konuşmaya başlayamama, telaffuz hataları, konuşurken ses yüksekliğini ayarlayamama gibi sosyal hayatı etkileyebilen zorlukları içermektedir [4].

Bir kişinin PH olup olmadığı basit bir test ile anlaşılabilir değildir. Nöroloji uzmanı bir doktor hastalığa teşhis koyabilmek ve rahatsızlığın başka bir hastalık durumundan kaynaklanıp kaynaklanmadığını anlamak amacıyla hastalardan biyokimyasal testler ve beyin tomografisi ister. Ek olarak, bacak ve kolların işlevsel yeterliliği, kas durumu, serbest yürüyüş ve dengeyi sağlayabilme durumlarını değerlendirmek için bazı fiziksel testler istemektedir [5]. Hastalar genellikle 60 ve üzeri yaşlarda [5] olduğu için istenen testler, bu yaşlardaki insanlara zor gelmektedir. Tüm bu zorluklar sebebiyle PH'ın tanısında daha basit ve güvenilir yöntemlere ihtiyaç duyulmaktadır [6-8].

PH'ın güvenilir ve hızlı bir şekilde teşhisinin yapılabilmesi için son yıllarda pek çok bilimsel çalışma yapılmıştır [6, 7, 9-20]. Bu çalışmaların en önemli amacı hem hastalar için fiziki zorlukları ortadan kaldırmak hem de klinik çalışanları üzerindeki iş yükünü azaltmaktır.

PH'ın erken evrelerinde ve en sık görülen rahatsızlıklarından biri vokal (konuşma ve ses) problemlerdir [9-11, 21]. Yakın geçmişte bireylerin konuşma kayıtları (sesleri) kaydedilerek PH teşhisi üzerinde çalışmalar gerçekleştirilmiştir [6-20]. Vokal çalışmalar literatürde belirgin bir yere sahiptir. Bu çalışmalarda genel olarak iki farklı veri seti kullanılmıştır. Bunlarda ilki 23 PH bireyden ve 8 sağlıklı bireyden 195 ses ölçümü ile ikinci veri seti ise 20 PH bireyden ve 20 sağlıklı bireyden alınan çoklu konuşma kayıtlarından oluşan halka açık veri seti kullanılmaktadır [11-18]. Diğer yeni

çalışmalarda genellikle küçük veri setleri kullanılarak PH teşhis edilmeye çalışılmıştır [6, 7, 9, 10, 19, 20]. Bunların dışında küçük olmayan veri setlerinde de yapılan çalışmalar vardır [8]. Ayarlanabilir Q-Faktör dalgacık dönüşümünün kullanıldığı veri seti ile bu dönüşümün kullanılmadığı veri setinde kıyaslamaya gidilmiştir. Dönüşümün düşük oranda da olsa performans kriterlerinden doğruluk oranını arttırdığı belirtilmiştir. Bu çalışmalarda veriler PH ve sağlıklı bakımdan dengesiz dağılıma sahip olmasına rağmen herhangi bir alt ya da üst örnekleme yapılmadan araştırma gerçekleştirilmiştir [6-20].

2018-2019 yıllarında PH teşhisi için makine öğrenmesi tabanlı birçok yeni sistemler geliştirilmiştir [6, 10, 20, 22]. Bunlardan en güncel yaklaşım olarak derin öğrenme tabanlı geliştirilen PH şiddetini tahmin eden modeldir [20]. Tensorflow (Makine öğrenme kütüphanesi) tabanında geliştirilen bu sistem, açık veri seti kaynağı olan UCI Machine Learning Repository'den alınan verileri kullanarak %62-81 doğruluk oranı ile çalışmaktadır [20]. 2019 yılında Sadek ve arkadaşlarının yapay sinir ağı tabanlı yaptığı bir çalışmada %93 başarı oranı ile PH tespit edilebilmiştir [10]. Çalışmada 31 bireyden alınan 195 ölçüm kullanılmıştır.

B.Karan ve arkadaşları ise 45 kişiden alınan 150 ölçüm ile destek vektör makineleri ve rastgele orman algoritması ile model kurmuşlardır. Küçük sayılabilecek bu veri setleri ile kurulan modelden %100 doğruluk oranı elde etmişlerdir [6]. Fakat eğitim ve test verileri birlikte test edildiğinde makine ezberleme yönüne gitmekte ve taraflı karar vermektedir. Veri seti eğitim / test setlerine uygun şekilde bölüdüğü zaman, önerilen modelin doğruluğu etkili bir biçimde azalmaktadır [6, 10, 11, 13, 17, 18]. Bu gibi küçük veri setlerinden alınan yüksek doğruluk oranları, büyük veri setlerinden alınamayacaktır. Daha uygun sonuç almak için daha büyük veri setinde birden fazla veri grupları kullanılması ve veri setinde dengeleme işlemi yapılması gerekmektedir [23-27].

C.Yücelbaş ve arkadaşları PH teşhisi için deneklerin yürüme verilerini kullanmıştır. Veri setinde yaş faktörüne göre gruplama yapılmıştır. Tavsiye edilen model Çift Yoğunluklu 1-D Dalgacık Dönüşümü yöntemi ile kurulmuştur [28]. H.Badem ve C.Yücelbaş 2019 yılında yapmış oldukları farklı çalışmalarda [29,30] yüksek doğruluk oranları bulmuşlardır. Çalışmalarda kapsamlı bir veri seti kullanılmıştır fakat dengeleme işlemi yapılmamıştır.

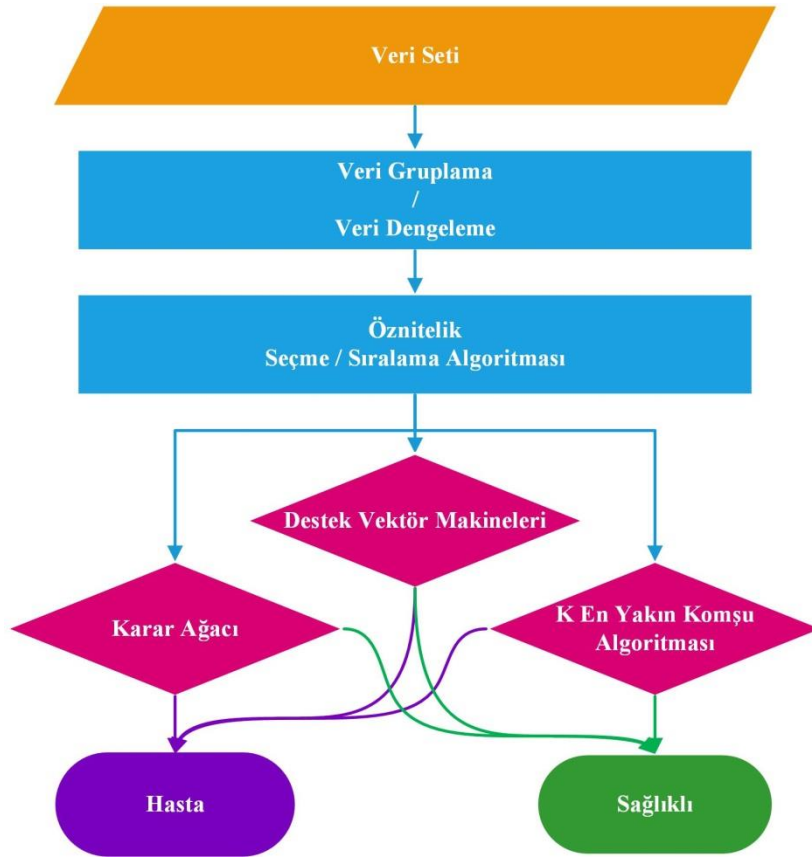
Literatürde genel olarak kullanılan veri setleri, makine öğrenmesinde kullanılmak üzere çıkarılan özellikler olarak vokal temel frekans, frekanstaki değişkenlik miktarları, genlikteki değişkenlik miktarları, gürültü ile ses tonu arasındaki bileşenlerin oranı, doğrusal olmayan dinamik değerler ve doğrusal olmayan temel frekans değerleri gibi benzer özellikleri içermektedir [11, 13, 14]. Bunlar gibi küçük veri setleri, PH teşhisinde hasta bireyler ile sağlıklı bireyleri ayırt etmede çok yüksek doğruluk oranları (%98-99) içermektedir [10, 11, 13, 17, 18]. H. Gürüler kompleks değerli yapay sinir ağı ve k ortalama algoritmasını birleştirerek oluşturduğu modelde %99 doğruluk değeri elde etmiştir [11]. Başka bir çalışmada regresyon, karar ağacı algoritması ve yapay sinir ağı ile kurulan modelde %98 doğruluk oranına ulaşılmıştır [18].

Bu çalışmanın amacı PH teşhisi için makine öğrenmesi tabanlı yeni yaklaşımlar sunmaktır. Halihazırda PH teşhisinde kullanılan yöntemler hastalar için fazladan efor sarf etmelerine ve klinik çalışanları için zaman kaybına sebep olabilmektedir. Bu makalede sunulan yöntem ile bireylerden alınan ses verileri kullanılarak PH'ın çok daha hızlı ve daha kolay bir şekilde teşhis edilmesi amaçlanmıştır. Bireylerden alınan ses verileri, makine öğrenme algoritmalarının kullanımıyla PH teşhisinde önemli rol oynamaktadır. Bu yöntemin uygulamasında hasta ve sağlıklı bireylerden alınan ses verileri kullanılarak makine öğrenmesi işlemi yapılıyor. Öğretilen ağ sayesinde yeni bireyin ses verilerinden hastalık teşhisine (varsa) gidilebiliyor. Bu amaçla, 252 bireyden elde edilen 756 ölçüm ile

752 öznitelik içeren veri seti kullanılmıştır. Bu veri seti literatürdeki diğer çalışmalarda kullanılan veri setlerinden daha kapsamlı bir veri setidir. Daha kapsamlı bir veri setinden alınan başarılı performans, daha etkili bir model oluşturmada yardımcı olmaktadır.

## II. MATERYAL VE YÖNTEM

Çalışma için şekil 1'deki akış şeması takip edildi. Veri setinde dengeleme işlemi yapıp ilgili öznitelik gruplarına ayrıldı. Oluşturulan veri grupları öznitelik seçme algoritması ile en ilgiliden en ilgisize doğru sıralandı. Bu sıralı veri grupları belli yüzdelik oranda öznitelik gruplarına bölündü ve her bir veri grubunun sınıflandırma algoritmaları ile performansları değerlendirildi. Veri setindeki düzenlemelerden performans değerlendirme aşamasına kadar gerekli işlemler Matlab programında gerçekleştirildi.



Şekil 1. Akış Şeması

### A. PARKİNSON HASTALIĞI VERİ SETİ

Çalışmamızda kullandığımız veri seti İstanbul Üniversitesi, Cerrahpaşa Tıp Fakültesi kaynaklı - Machine Learning Repository- (UCI)'den alınmıştır. Veri seti 188 (107 Erkek ve 81 Kadın) Parkinson hastası bireyler ile 64 (23 Erkek ve 41 Kadın) kontrol grubu olan Parkinson hastası olmayan

bireylerden, toplam 252 bireyden oluşmaktadır. Bireylerin yaşları 33 ile 87 arasında değişmektedir. Veri setinde etiket 1 ve 0'lerden oluşmaktadır (1 hasta grubu, 0 sağlıklı grubu). 252 bireyin her birinden 3 defa /a/ sesli harfini söylemesi istenmiştir ve toplamda 756 ölçüm elde edilmiştir. Bu elde edilen kayıtlardan bir tanesi etiket olmak üzere 753 öznitelik oluşturulmuştur.

## B. VERİ ÖNİŞLEME

Literatürde, veri setini analize hazır hale getirmek için Han ve Kamber (2006) tarafından oluşturulan bazı adımlar mevcuttur [31]. Bu çalışmada uygulanan veri önışleme adımları aşağıda belirtilmiştir.

### B.1. Veri Setini İlgili Öznitelik Gruplarına Ayırma

Ham veri setindeki 752 öznitelik, bazı öznitelik gruplarından elde edildi. Bunlar Baseline Features (Temel Özellikler), Intensity Parameters (Yoğunluk Parametreleri), Formant Frequencies (Fonetik Biçimsel Frekanslar), Bandwidth Parameters (Bant Genişliği Parametreleri), Vocal Fold (Ses Kıvrım Parametreleri), MFCC (Mel-Frekanslı Cepstral Katsayıları) ve Wavelet Features (Dalgacık Özellikleri) dir. Benzer öznitelikler bulundurduğu için 3 grup, Time Features (Zamansal Özellikler) olarak birleştirildi (Intensity Parameters, Formant Frequencies ve Bandwidth Parameters). Böylelikle işleme alınacak 5 temel öznitelik grubu ve tüm özniteliklerden oluşan bir grup olmak üzere toplamda 6 veri grubu Tablo 1'de gösterildiği gibi oluşturulmuştur.

*Tablo 1. Veri Seti Öznitelik Grupları*

Veri Grubu İsmi	Öznitelik Sayısı
Baseline	21
Mfcc	84
Time	11
Vocal	22
Wavelet	614
Tüm	752

### B.2. Veri Setini Dengeli Hale Getirme

Veri setinde bulunan etiket sınıfı değerleri arasında eşitlik durumu olmaması halinde veri seti için dengesiz veri seti denilmektedir [23]. Yapılan çalışmada kullanılan veri seti dengesiz ise, doğruluk değerleri performans değerlendirmede yanıltıcı karar vermeye sebep olabilmektedir [23]. Bu olumsuz durumu izale edebilmek için sistematik örnekleme metodu kullanılmıştır [24]. Bu metot içerisindeki alt-örnekleme (Undersampling) yöntemi ile dengeleme işlemi yapıldı ve dengesiz durumdan kurtarıldı. Alt-örnekleme yönteminde, sayıca fazla olan etiket sınıfı, sayıca az olan etiket sınıfına eşitlenir [23]. Bu çalışmadaki veri setinde 756 ölçümden, 192'si sağlıklı etikete (0) ve 564'ü hasta etikete (1) sahiptir. Bu yüzden dengeleme sonucunda 192 sağlıklı ve 192 hasta etiketi elde edilmiştir.

*Tablo 2. Veri Seti Dengeleme Tablosu*

Veri Seti	Bireyler	Ölçüm Sayısı
Ham Veri Seti	Hasta	564
	Sağlıklı	192
Dengeli Veri Seti	Hasta	192
	Sağlıklı	192

## C. ÖZELLİK SEÇME / SIRALAMA ALGORİTMASI

### C.1. Fisher Öznitelik Sıralama Algoritması

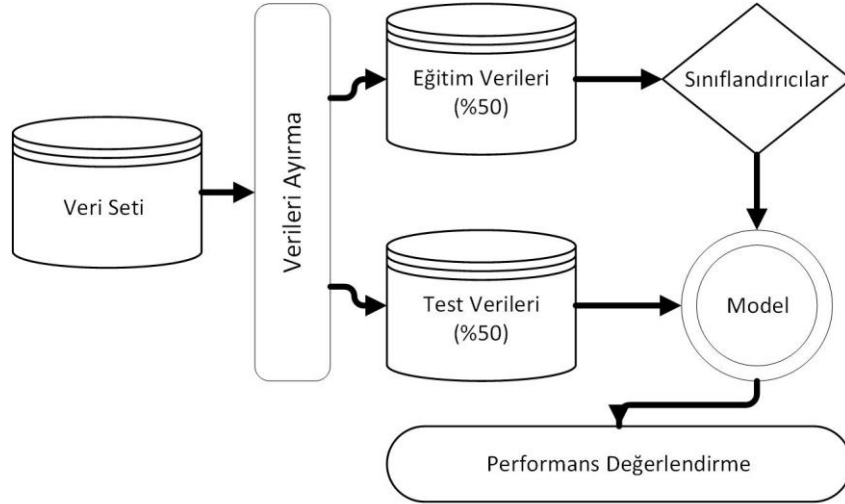
Öznitelik sayısının makine öğrenmesi performansında olumlu ve olumsuz etkileri mevcuttur. Olumsuz etkileri izale etmek için öznitelik seçme işlemine gidilir. Bu işlem ile herhangi bir özneliğin etiket tahminindeki etki gücüne göre ilgiliden ilgisize doğru bir sıralama yapılır. Araştırmacı, bu en ilgiliden en ilgisize doğru sıralanmış veri setindeki özneliklerin istediği kadarını çalışmasına dahil edebilir. Böylece gereksiz verileri kullanmayıp daha doğru sonuç alabilir ve daha hızlı bir program döngüsüne sahip olabilir. Bu çalışmada öznitelik seçim algoritması kullanıldı ve Tablo 3’de gösterildiği gibi en ilgiliden başlayacak şekilde veri setinin %5,%10. . . %50’si alınarak performans değerlerine göre model kurulumu gerçekleştirildi.

*Tablo 3. Veri Gruplarında Öznitelik Sıralaması*

Veri Grupları	%5	%10	%15	%20	%25
Baseline	1	2	3	4	6
MFCC	4	9	13	17	21
Time	1	1	2	2	3
Vocal	1	2	3	5	6
Wavelet	31	62	92	123	154
Tüm	38	75	113	151	188
Veri Grupları	%30	%35	%40	%45	%50
Baseline	7	8	9	10	11
MFCC	26	30	34	38	43
Time	4	4	5	5	6
Vocal	7	8	9	10	12
Wavelet	185	215	246	277	308

## D. SINIFLANDIRMA ALGORİTMALARI

Çalışmamızdaki sınıflandırma işlemleri Karar ağacı (DT) , Destek vektör makineleri (DVM) ve K en yakın komşu (kNN) algoritmalarıyla gerçekleştirilmiştir. Sınıflandırma işlemi için şekil 2'deki akış şeması adımları uygulanmıştır. Veri setinde sınıflandırma işlemi yapmak için verilerin yarısı model oluşturmada diğer yarısı modeli test etme aşamasında kullanılmıştır. Her bir veri grubu için sistematik örnekleme metodu yardımı ile eğitim veri seti oluşturulmuştur. Geriye kalan veri grubu test aşamasında kullanılmıştır. Test verileri üzerinden kurulan modelin performans değerlendirme kriterleri sınımlanmıştır.



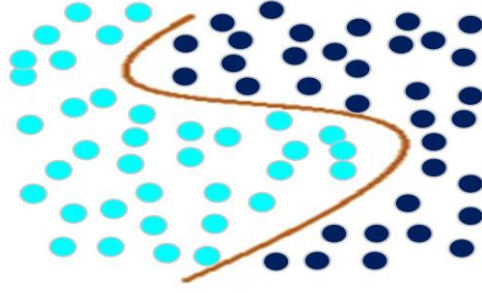
Şekil 2. Veri Setini Eğitim ve Test Veri Gruplarına Ayırma

### D.1. Karar Ağacı

Karar ağacı (DT) algoritmasının temel yapısında kök, dal, düğüm ve yapraklar bulunur. Ağaç yapısı oluşturulurken her bir öznelik bir düğüm ile ilişkilendirilir. Kök ve düğümler arasında dallar bulunur. Her bir düğümden dallar aracılığı ile diğer düğüme geçilir. Ağaçta karar, gidilen nihai yaprağa göre verilir [32]. Karar ağacı yapısı oluşturulmasında temel mantık, ulaşılan her bir düğümden ilgili soruları sorup, verilen cevaplara göre dallar üzerinden nihai yaprağa en kısa yol ve süre içerisinde ulaşmak olarak özetlenebilir. Böylece sorulardan elde edilen cevaplara göre karar kuralları / model oluşturur. Eğitilen bu ağaç yapısı farklı test verileri ile doğruluğu sınanır ve uygun sonuç üretirse model kullanılır. Karar ağacı modeli şekil 3'te gösterilmiştir.

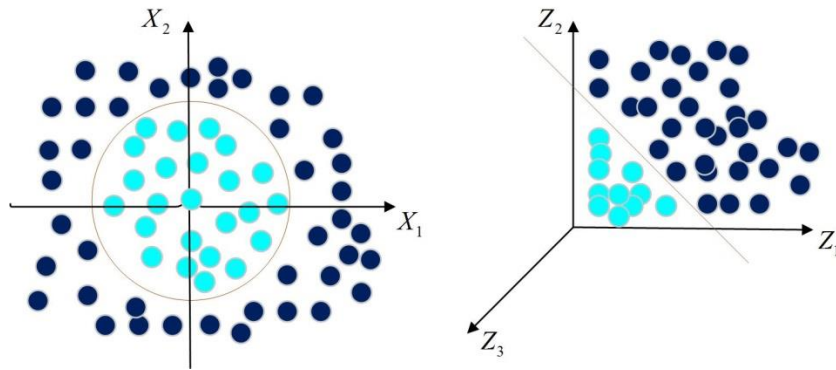






Şekil 5. Doğrusal Ayrılmayan Nesnelere

Şekil 5'teki gibi nesnelere sahip bir veri setinde sınıflandırma yapmak için, boyut dönüşümü yapılmaktadır.



Şekil 6. İki Boyutlu Nesnelere Üç Boyuta Aktarımı

Şekil 6'da iki boyutlu haritadan üç boyutlu harita üzerine geçirilmiş nesnelere görünmektedir. Yeni durumda fonksiyon uzayı Denklem 1'deki eşitliğe ulaşır.

$$\mathcal{G}: R^2 \rightarrow R^3 : (X_1, X_2) \rightarrow (Z_1, Z_2, Z_3) := (X_1^2, \sqrt{2X_1X_2}, X_2^2) \quad (1)$$

Fonksiyondaki girdi vektörlerinin (X'ler) iç çarpımları alındığında Denklem 2 elde edilir.

$$= (X_1^2 X_2'^2 + 2X_1 X_1' X_2 X_2' + X_2^2 X_2'^2)^2 \quad (2)$$

$$f(X) = \text{sign}((W^*, X_1) + b^*) \quad (3)$$

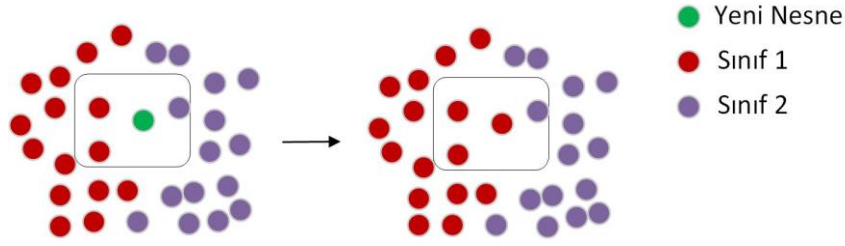
$$f(X) = \text{sign}\left(\sum_{i=1}^{\infty} \gamma_i \alpha_i (\mathcal{G}(X_i), \mathcal{G}(X_j))\right) \quad (4)$$

Denklem 3'te bulunan w normal düzlemini ifade ederken, b sabit sayıyı ve Denklem 4'deki  $\alpha$  Lagrange çarpanını,  $\gamma$  ise nesnelere bulunduğu sınıfı temsil etmektedir [35]. Denklem 3 ve 4'teki sign ifadesi Signum (İşaret) fonksiyonudur. Bu fonksiyon içerisindeki ifadeye göre -1,0 veya 1

değerini alır. Fonksiyon, içerisindeki değer negatif ise -1, değer pozitif ise 1 ve içerisindeki değer sıfıra eşit olduğu durumda da 0 değerini alır. Fonksiyonun değerine göre nesne konumu tayin edilir. Böylece nesnelere iki boyutlu haritadan üç boyutlu haritaya taşınmış olur. Yeni sınıflandırılacak nesnelere Denklem 4'e göre sınıflandırılır.

### D.3. K En Yakın Komşu Algoritması

Sınıflandırma işleminde kullanılan K-En Yakın Komşu (kNN) Algoritması oldukça basit bir mantık üzerine temellendirilmiştir. Veri setinde bulunan her bir nesne hangi komşusu ile arasındaki mesafesi en az ise o komşunun sınıfı ile sınıflandırılır. Her bir nesne için komşuluk mesafeleri hesaplanarak K-En Yakın Komşu Algoritması oluşturulmuş olur. Algoritmada ne kadar komşuya bakılarak sınıflandırma yapılacağını K parametresi ifade eder [36]. Veri setinde bulunan her bir nesne için etrafındaki K tane komşusunun hangi sınıfa ait olduğuna bakılır. Şekil 7'de belirtildiği gibi komşuları en çok hangi sınıfta ise nesne de o sınıfa dahil edilir. Eşitlik durumu oluşmaması için K değeri genel olarak tek sayılardan seçilir [37]. Bu çalışmada K değeri 3 olarak seçilmiştir.



Şekil 7. K En Yakın Komşu Algoritması

Nesneler arasındaki mesafelerin belirlenmesinde Öklid (Denklem 5), manhattan (Denklem 6) yada minkowski (Denklem 7) gibi uzaklık fonksiyonları kullanılır [38]. Sınıflandırma algoritmalarında bu fonksiyonlardan en sık kullanılanı Öklid uzaklık fonksiyonudur [39]. Öklid, nesnelere arasındaki doğrusal uzaklığı ölçer. Nesnelere birisi  $Q = (X_1, X_2, X_3 \dots X_n)$ , diğeri  $S = (Y_1, Y_2, Y_3 \dots Y_n)$  ise bu iki nesne arasındaki uzaklık Denklem 5'teki formül ile hesaplanmaktadır.

$$\sqrt{\sum_{i=1}^n (X_i - Y_i)^2} \quad (5)$$

Manhattan uzaklık fonksiyonu nesnelere arasındaki konum farklarının mutlak değerlerinin toplamıdır [39]. Yine  $Q = (X_1, X_2, X_3 \dots X_n)$  ve  $S = (Y_1, Y_2, Y_3 \dots Y_n)$  nesnelere manhattan uzaklık fonksiyonu ile hesaplanması Denklem 6'da gösterilmiştir.

$$\sum_{i=1}^n |X_i - Y_i| \quad (6)$$

Minkowski uzaklık fonksiyonu Öklid uzayında tanımlanmıştır. Bu fonksiyon Öklid ve manhattan uzaklık fonksiyonlarının genelleştirilmiş halidir [39]. İki nesne olan  $Q = (X_1, X_2, X_3 \dots X_n)$  ve

$S = (Y_1, Y_2, Y_3 \dots Y_n)$  nesnelarının minkowski uzaklık hesaplaması Denklem 7’de gösterilmiştir. Minkowski uzaklık fonksiyonundaki q parametresi diđer uzaklık fonksiyonlarını tanımlamak için kullanılmaktadır. Minkowski uzaklık fonksiyonu q deęeri 1 seęilirse Manhattan’a, q deęeri 2 seęilirse Öklid’e benzeyecektir.

$$\left( \sum_{i=1}^n (|X_i - Y_i|)^q \right)^{1/q} \quad (7)$$

## E. PERFORMANS DEęERLENDİRME KRİTERLERİ

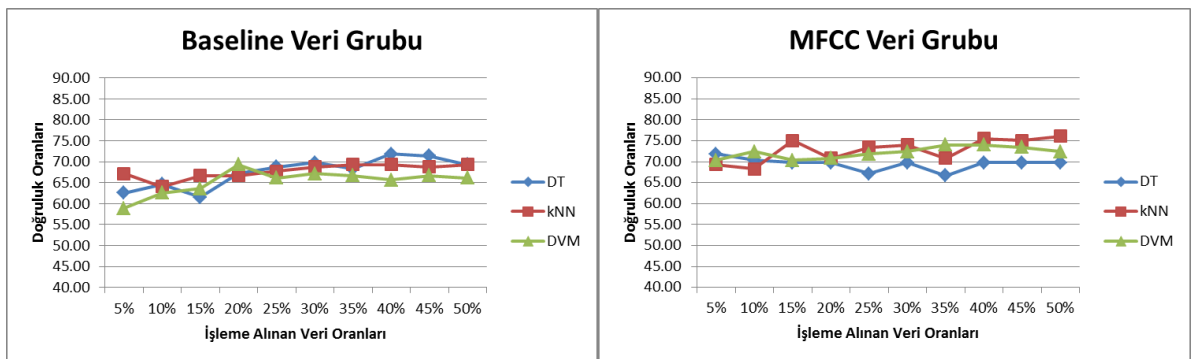
Çalışmamızda kurulan her bir model için Doğruluk, Duyarlılık, Özgünlük, F-Ölçümü, Kappa Katsayısı, AUC (ROC eğrisinin altında kalan alan) gibi performans deęerlendirme kriterleri incelenmiştir. Bu performans deęerlendirme kriterleri kullanılan üç sınıflandırıcı (Karar Ağacı, Destek Vektör Makineleri, K En Yakın Komşu Algoritması) çıktılarına uygulanmıştır. Veri seti sınıflandırılırken eğitim-test oranları %50-%50 şeklinde belirlenmiştir (Tablo 4).

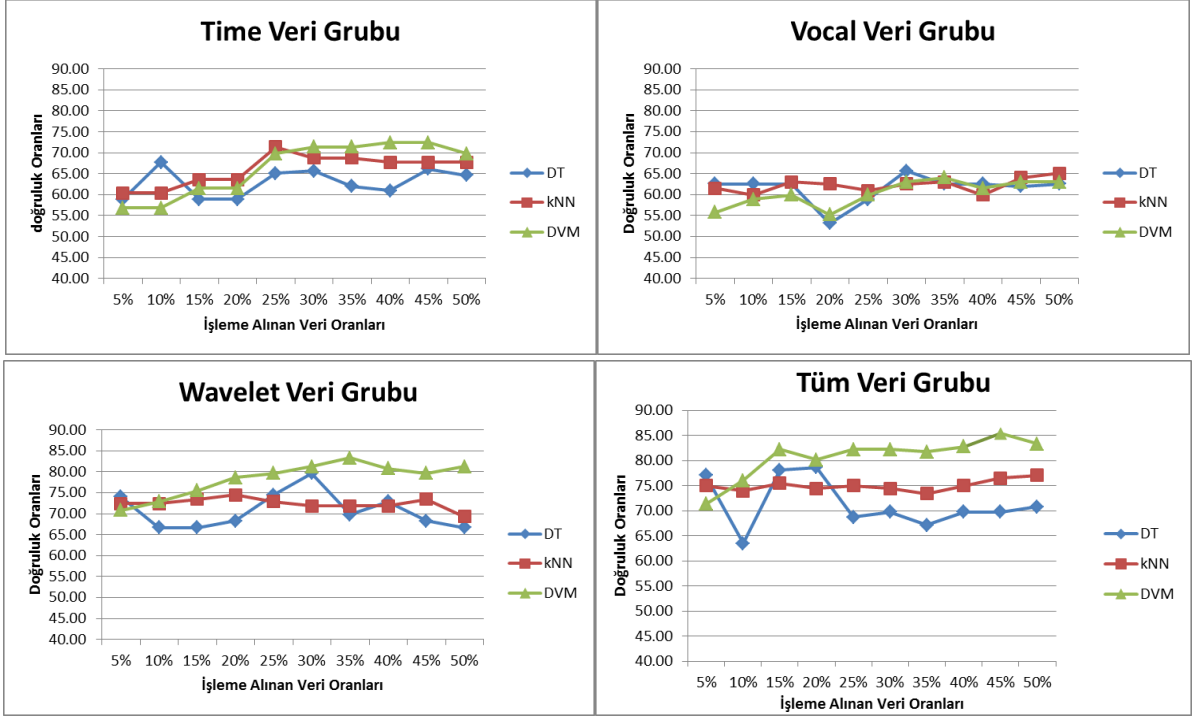
*Tablo 4. Eğitim ve Test Seti Dağılımı*

	Eğitim (%50)	Test (%50)	Toplam
<b>Hasta</b>	86	86	192
<b>Saęlıklı</b>	86	86	192
<b>Toplam</b>	192	192	384

## III. BULGULAR

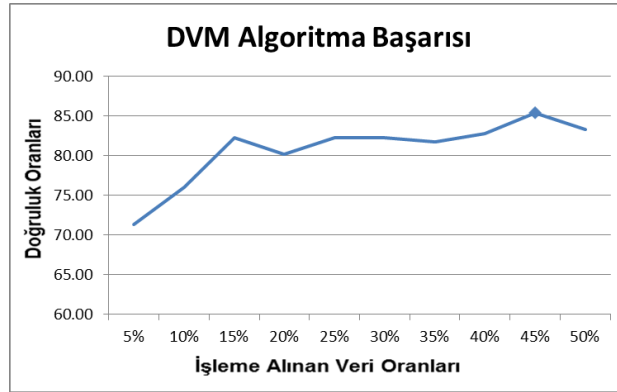
Çalışmamızda PH’in makine öğrenmesi ile teşhis edilebilmesi hedeflenmiştir. Bu amaç için sınıflandırıcı algoritmaları kullanılmıştır. Sınıflandırıcı algoritmalar belirli veri gruplarına uygulanmıştır ve uygun performans deęerleri elde edilmiştir.





Şekil 8. Farklı Veri Grupları İçin Farklı Sınıflandırıcı Başarıları

Şekil 8 ve Tablo 5'te sınıflandırıcı başarıları görülmektedir. Veri gruplarının en yüksek doğruluk oranını veren sınıflandırıcıları sonuçları; Baseline için %71, MFCC için %75, Time için %72, Vocal için %64, Wavelet için %80, Tüm için %85 oranındadır. Her veri grubu için belli bir doğruluk oranında sınıflandırıcı başarıları mevcuttur. Fakat en yüksek doğruluk oranı, her veri grubunu içeren Tüm adlı veri grubunda görülmüştür. Sınıflandırma algoritmalarından en yüksek doğruluk oranı Destek Vektör Makineleri'nde (DVM) görülmüştür.



Şekil 9. Tüm Adlı Veri Grubunda DVM Doğruluk Oranları

En yüksek doğruluk oranı (%85) , Tüm adlı gruptaki verilerin %45'i alındığında ve sınıflandırma için Destek Vektör Makineleri algoritması kullanıldığında elde edilmiştir. Veri setinin %50 değil de %45'i eğitime alındığında daha iyi sonuç vermiş olması, veri setinde en ilgiliden en ilgisize doğru öznitelik sıralaması yapılmasından kaynaklandığı söylenebilir. Diğer performans kriterlerine bakıldığı zaman Duyarlılık, F-ölçümü, Kappa ve Auc için Destek Vektör Makinesi en yüksek başarıyı sağladığı görülür. Sadece Özgünlük kriteri için Karar Ağacı algoritması az bir farkla (0.02) daha yüksek başarı

sağlamıştır. Başarılı sonuç alınan grubun diğer performans kriterleri; Duyarlılık: 0.94, Özgünlük: 0.77, F-Ölçümü: 0.85, Kappa: 0.71, Auc: 0.85 olarak hesaplanmıştır.

## **V. SONUC**

Çalışmamızda PH'ı teşhis etmede makine öğrenmesinden faydalanmak amaçlanmıştır. Makine öğrenmesinde kullanılan veri seti sadece hastaların ses kayıtlarının analizlerinden oluşmuştur. Bu sayede teşhis süreci hem daha kısa hem daha az maliyette olacaktır. Ayrıca klinik çalışanlarının iş yükünü azaltacak ve hastaların daha kolay bir teşhis süreci geçirmelerini sağlayacaktır.

Literatürde PH teşhisi için birçok çalışma mevcuttur [6-20, 26]. C.Yücelbaş ve arkadaşları PH teşhisi için denekleri yürüterek veri seti oluşturmuştur. Oluşturulan veri seti, yaş faktörüne göre gruplama yapılarak Çift Yoğunluklu 1-D Dalgacık Dönüşümü yöntemiyle analiz edilmiştir [28]. Fakat küçük veri setinde ve az özniteliklerle çalışılmıştır.

Bir diğer çalışmada H.Badem ve arkadaşları Yapay Sinir Ağlarını kullanarak %87 doğruluk oranı elde edilen bir model kurmuşlardır. Kurulan modelde iki veri seti kullanılmıştır. Veri setlerinden biri 23 öznitelikten oluşurken diğeri 26 öznitelikten oluşmaktadır [27].Kurdukları yüksek doğruluk oranına sahip modeller büyük veri setlerinde aynı oranda doğruluk sağlayamayacaktır. Kullandığımız veri setinde veri sayısı oldukça büyüktür. Bu yüzden bu makalede oluşturulan modeller daha güvenilir sonuçlar üretebilir.

Literatürde hali hazırda bulunan pek çok veri seti dengesiz dağılıma sahiptir. Yapılan çalışmalarda bu dengesizlik giderilmeden model oluşturulmuştur [6-11, 16-20]. 2019 yılında H.Badem [29] ve C.Yücelbaş [30] yapmış oldukları farklı çalışmalarda yüksek doğruluk oranları bulmuşlardır. Çalışmalarda kapsamlı bir veri seti kullanılmıştır fakat dengeleme işlemi yapılmamıştır. Bu çalışmada kullanılan veriler alt örnekleme yöntemi ile dengelenerek model oluşturulduğu için daha kararlı çalışacağını düşünmekteyiz. Dengesiz verilerle oluşturulan modellerde sistem miktar bakımından fazla olan verilere yatkın sonuçlar üretir [23-25]. Bu çalışmada önerilen modeller dengeli veri setleri ile kurulduğu için literatüre göre bir adım öndedir.

Literatürdeki bazı çalışmalarda, verilen sonuçlar eğitim ve test performanslarının ortalaması olarak verilmektedir [13, 40]. Bu bakımdan, bu çalışmaların sonuçları tartışılmalıdır. Bu makalede ise test verileri bağımsız olarak değerlendirilmiştir. Bu yüzden doğruluk oranları yüksek olan makalelere [13, 40] nispeten biraz düşük çıkmıştır. Ancak kabul edilir değerlerdir ve diğer çalışmalara göre daha güvenilirdir [13, 40].

Bu çalışmada öznitelik seçme algoritmasına göre sıralanmış özniteliklerin ilk %45'i alındığında en iyi sonuca ulaşılmıştır. Daha fazla veri grubu ile çalışıldığında hem doğruluk oranları düşmekte hem de döngü hızı yavaşlamaktadır. Her bir veri grubu için sınıflandırma algoritmalarına bakıldığında en iyi performans değerleri Destek Vektör Makineleri algoritmasında görülmüştür (Tablo 5). Tüm adlı gruptaki veri setinin %45 lik en ilgili öznitelikleri ile sınıflandırılma işlemi yapıldığında bu sonuçlar elde edilmiştir (Doğruluk:%85, Duyarlılık: 0.94, Özgünlük: 0.77, F-Ölçümü: 0.85, Kappa: 0.71, Auc: 0.85). Oluşturulan model ile PH teşhisindeki zorlu ve maliyetli olan teşhis işlemi kolaylaştırarak doktora tıbbi karar destek olabileceği sonucuna ulaşılmaktadır.

*Tablo 5. Her Veri Grubu İçin Sınıflandırıcı Başarı Tablosu*

Data seti	Baseline Veri Grubu			MFCC Veri Grubu		
Performans Kriterleri	Sınıflandırıcı Algoritmaları					
	DT	kNN	DVM	DT	kNN	DVM
Doğruluk Oranı (%)	71.35	68.75	66.67	69.79	75.00	73.44
Duyarlılık	0.75	0.70	0.72	0.60	0.71	0.70
Özgüllük	0.68	0.68	0.61	0.79	0.79	0.77
F-Ölçümü	0.71	0.69	0.66	0.69	0.75	0.73
Kappa	0.43	0.38	0.33	0.40	0.50	0.47
AUC	0.71	0.69	0.67	0.70	0.75	0.73
Data seti	Time Veri Grubu			Vocal Veri Grubu		
Performans Kriterleri	Sınıflandırıcı Algoritmaları					
	DT	kNN	DVM	DT	kNN	DVM
Doğruluk Oranı (%)	66.15	67.71	72.40	61.98	64.06	63.02
Duyarlılık	0.73	0.82	0.80	0.64	0.58	0.64
Özgüllük	0.59	0.53	0.65	0.60	0.70	0.63
F-Ölçümü	0.65	0.65	0.72	0.62	0.64	0.63
Kappa	0.32	0.35	0.45	0.24	0.28	0.26
AUC	0.66	0.68	0.72	0.62	0.64	0.63
Data seti	Wavelet Veri Grubu			Tüm Veri Grubu		
Performans Kriterleri	Sınıflandırıcı Algoritmaları					
	DT	kNN	DVM	DT	kNN	DVM
Doğruluk Oranı (%)	68.23	73.44	79.69	69.79	76.56	<b>85.42</b>
Duyarlılık	0.81	0.76	0.97	0.60	0.81	0.94
Özgüllük	0.55	0.71	0.63	0.79	0.72	0.77
F-Ölçümü	0.66	0.73	0.76	0.69	0.76	<b>0.85</b>
Kappa	0.36	0.47	0.59	0.40	0.53	<b>0.71</b>
AUC	0.68	0.73	0.80	0.70	0.77	<b>0.85</b>

## **V. KAYNAKLAR**

- [1] A. Wood-Kaczmar, S. Gandhi, and N.W. Wood, "Understanding the molecular causes of Parkinson's disease," Trends in Molecular Medicine, vol 12, no 11, pp 521–528, 2006.
- [2] Roger A Barker and Stephen B. Dunnett, "Functional integration of neural grafts in Parkinson's disease," 1999.
- [3] William M. McDonald, Paul E. Holtzheimer, and Eve H. Byrd, "The diagnosis and treatment of depression in parkinson's disease," Current Treatment Options in Neurology, vol 8, no 3, pp 245–255, 2006.
- [4] Pratibha Surathi, Ketan Jhunhunwala, Ravi Yadav, and PramodKumar Pal, "Research in Parkinson's 7 disease in India: A review," Annals of Indian Academy, 19-1, 2016.

- [5] Parkinsondernegi, “Doktorunuz Parkinson Hastalığı Tanısını Nasıl Koyar?”, Erişim: 27.01.2020, <http://parkinsondernegi.com/>
- [6] Biswajit Karan, Sitanshu Sekhar Sahu, and Kartik Mahto, “Parkinson disease prediction using intrinsic mode function based features from speech signal,” *Biocybernetics and Biomedical Engineering*, may 2019.
- [7] Richa Mathur, Vibhakar Pathak, and Devesh Bandil, “Parkinson Disease Prediction Using Machine Learning Algorithm” pp 357–363, 2019.
- [8] C. Okan Sakar, Gorkem Serbes, Aysegul Gunduz, Hunkar C. Tunc, Hatice Nizam, Betul Erdogan Sakar, Melih Tutuncu, Tarkan Aydin, M. Erdem Isenkul, and Hulya Apaydin, “A comparative analysis of speech signal processing algorithms for Parkinson’s disease classification and the use of the tunable Q-factor wavelet Transform,” *Applied Soft Computing*, vol 74, pp 255–263, 2019.
- [9] Jing Tang, Bao Yang, Matthew P. Adams, Nikolay N. Shenkov, Ivan S. Klyuzhin, Sima Fotouhi, Esmail Davoodi-Bojd, Lijun Lu, Hamid Soltanian-Zadeh, Vesna Sossi, and Arman Rahmim, “Artificial Neural Network Based Prediction of Outcome in Parkinson’s Disease Patients Using DaTscan SPECT Imaging Features,” *Molecular Imaging and Biology*, mar 2019.
- [10] Ramzi M. Sadek, Salah A. Mohammed, Abdul Rahman K. Abunbehan, Abdul Karim H. Abdul Ghattas, Majed R. Badawi, Mohamed N. Mortaja, Bassem S. Abu-Nasser, and Samy S. Abu-Naser, “Parkinson’s Disease Prediction Using Artificial Neural Network,” 2019.
- [11] Hüseyin Gürüler, “A novel diagnosis system for Parkinson’s disease using complex-valued artificial neural network with k-means clustering feature weighting method,” *Neural Computing and Applications*, vol 28, no 7, pp 1657–1666, 2017.
- [12] Ömer Eskidere, “A Comparison of Feature Selection Methods for Diagnosis of Parkinson’s Disease from Vocal Measurements,” *Journal Engineering and Natural Science*, vol 30, pp 402–414, 2012.
- [13] C. Okan Sakar and Olcay Kursun, “Telediagnosis of Parkinson’s Disease Using Measurements of Dysphonia,” *Journal of Medical Systems*, vol 34, no 4, pp 591–599, 2010.
- [14] A. Tsanas, M. A. Little, P. E. McSharry, and L. O. Ramig, “Nonlinear speech analysis algorithms mapped to a standard metric achieve clinically useful quantification of average Parkinson’s disease symptom severity,” *Journal of The Royal Society Interface*, vol 8, no 59, pp 842–855, 2011.
- [15] Max A Little, Patrick E McSharry, Eric J Hunter, Jennifer Spielman, and Lorraine O Ramig, “Suitability of dysphonia measurements for telemonitoring of Parkinson’s disease,” *IEEE transactions on bio-medical engineering*, vol.56, no.4, pp.1015, 2009.
- [16] Musa Peker, “A decision support system to improve medical diagnosis using a combination of k-medoids clustering based attribute weighting and SVM,” *Journal of Medical Systems*, vol.40, no.5, pp. 116, 2016.

- [17] Betül Erdogdu Sakar, Gorkem Serbes, and C Okan Sakar, “Analyzing the effectiveness of vocal features in early telediagnosis of Parkinson’s disease” PloS one, vol.12, no.8, pp. 0182428, 2017.
- [18] Musa Peker, Baha Sen, and Dursun Delen, “Computer-Aided Diagnosis of Parkinson’s Disease Using Complex-Valued Neural Networks and mRMR Feature Selection Algorithm” Journal of healthcare engineering, vol.6, no. 3, pp.281–302, 2015.
- [19] Kyoungjune Pak, Heeyoung Kim, Ju Won Seok, Myung Jun Lee, Seunghyeon Shin, Keunyoung Kim, Jae Meen Lee, Youngduk Seo, Bum Soo Kim, Sungmin Jun, and In Joo Kim, “Prediction of future weight change with dopamine transporter in patients with Parkinson’s disease,” Journal of Neural Transmission, vol. 126, no.6, pp. 723–729, 2019.
- [20] Srishti Grover, Saloni Bhartia, Akshama, Abhilasha Yadav, and Seeja K.R, “Predicting Severity Of Parkinson’s Disease Using Deep Learning,” Procedia Computer Science, vol.132, pp. 1788–1794, 2018.
- [21] Sukru Torun “Parkinsonlularda Konuşma Fonksiyonunun Subjektif ve Objektif (Elektrolaringografik) Yöntemlerle incelenmesi,” 1991.
- [22] Timothy J. Wroge, Yasin Ozkanca, Cenk Demiroglu, Dong Si, David C. Atkins, and Reza Hosseini Ghomi, "Parkinson’s Disease Diagnosis Using Machine Learning and Voice,” In 2018 IEEE Signal Processing in Medicine and Biology Symposium (SPMB), pp 1–7. IEEE, dec 2018.
- [23] Elif Kartal and Zeki Ozen, “Dengesiz Veri Setlerinde Sınıflandırma,” pp. 109–131, 2017.
- [24] Reha Alpar, “Spor Sağlık Ve Eğitim Bilimlerinden Örneklerle UYGULAMALI İSTATİSTİK VE GEÇERLİK GÜVENİRLİK,” DETAY YAYINCILIK, 5 edition, 2018.
- [25] Salvador García, Sergio Ramírez-Gallego, Julián Luengo, José Manuel Benítez, and Francisco Herrera, “Big data preprocessing: methods and prospects,” Big Data Analytics, vol. 1, no. 1, pp. 9, 2016.
- [26] YÜCELBAŞ, Ş, YÜCELBAŞ, C, “Temel Bileşen Analizi Yöntemleri Kullanarak Parkinson Hastalığının Otomatik Teşhisi,” Avrupa Bilim ve Teknoloji Dergisi, c. 16, ss. 294-300, 2019.
- [27] H. Badem, A. Caliskan, A. Basturk and M. E. Yuksel, "Classification and diagnosis of the parkinson disease by stacked autoencoder," 2016 National Conference on Electrical Electronics and Biomedical Engineering (ELECO), pp. 499-502, Bursa, 2016.
- [28] YÜCELBAŞ, C, YÜCELBAŞ, Ş, “Çift Yoğunluklu 1-D Dalgacık Dönüşümü Kullanılarak Parkinson Hastalığının Yaş Faktörüne Göre Tespit Edilmesi,” Avrupa Bilim ve Teknoloji Dergisi, 2019.
- [29] H. Badem, “PARKİNSON HASTALIĞININ SES SİNYALLERİ ÜZERİNDEN MAKİNE ÖĞRENME Sİ TEKNİKLERİ İLE TANIMLANMASI,” Niğde Ömer Halisdemir Üniversitesi Mühendislik Bilimleri Dergisi, c.8, s.2, ss. 630-637, 2019.



- [30] YÜCELBAŞ, C, YÜCELBAŞ, Ş, “AQDD Özelliklerine BBA Yöntemleri Uygulanarak Parkinson Hastalığının Otomatik Teşhisi,” Bilecik Şeyh Edebali Üniversitesi Fen Bilimleri Dergisi, 2019.
- [31] Elif Kartal, “Sınıflandırmaya Dayalı Makine Öğrenmesi Teknikleri ve Kardiyolojik Risk Değerlendirmesine İlişkin Bir Uygulama,” PhD thesis, 2015.
- [32] J. R. (John Ross) Quinlan and J. Ross, “C4.5: programs for machine learning,” Morgan Kaufmann Publishers, 1993.
- [33] M.A. Hearst, S.T. Dumais, E. Osuna, J. Platt, and B. Scholkopf, “Support vector machines,” IEEE Intelligent Systems and their Applications, vol.13, no.4, pp18–28, 1998.
- [34] Sevgi AYHAN and Şenol ERDOĞMUŞ. Destek Vektör Makineleriyle Sınıflandırma Problemlerinin Çözümü İçin Çekirdek Fonksiyonu Seçimi. Eskişehir Osmangazi Üniversitesi İktisadi ve İdari Bilimler Dergisi, c. 9, s.1, ss. 75–201, 2014.
- [35] Soman Kp, R Loganathan, and Ajay Vadakkepatt, “Machine learning with SVM and other kernel methods,” 2009.
- [36] Pádraig Cunningham and Sarah Jane Delany, “k-Nearest Neighbour Classifiers,” Technical report, 2007.
- [37] BALABAN M. Erdal and KARTAL Elif, “Veri Madenciliği ve Makine Öğrenmesi Temel Algoritmaları ve R Dili ile Uygulamaları,” Çağlayan Kitabevi, 2 edition, 2018.
- [38] M Özgür Dolgun, Özdemir T. Güzel, and Oğuz Doruk, “Veri madenciliği’nde yapısal olmayan verinin analizi: Metin ve web madenciliği,” İstatistikçiler Dergisi, pp. 48-58, 2009.
- [39] E. Taşci, A. Onan, "K-en yakın komşu algoritması parametrelerinin sınıflandırma performansı üzerine etkisinin incelenmesi," Akademik Bilişim 2016, Aydın, 2016.
- [40] Betül Erdogdu Sakar, M. Erdem Isenkul, C. Okan Sakar, Ahmet Sertbas, Fikret Gurgun, Sakir Delil, Hulya Apaydin, and Olcay Kursun, “Collection and Analysis of a Parkinson Speech Dataset With Multiple Types of Sound Recordings,” IEEE Journal of Biomedical and Health Informatics, vol.17, no.4, pp. 828–834, 2013.