



T.C.

AKDENİZ ÜNİVERSİTESİ  
EĞİTİM BİLİMLERİ ENSTİTÜSÜ  
EĞİTİM BİLİMLERİ  
ANA BİLİM DALI

YÜKSEK  
LİSANS  
TEZİ

EĞİTİM ALANINDA MAKİNE  
ÖĞRENİMİ SINIFLANDIRMA  
ALGORİTMALARININ İNCELENMESİ

Canay CAN

EĞİTİMDE ÖLÇME VE DEĞERLENDİRME  
TEZLİ YÜKSEK LİSANS PROGRAMI

Antalya, 2021

**T.C.**  
**AKDENİZ ÜNİVERSİTESİ**  
**EĞİTİM BİLİMLERİ ENSTİTÜSÜ**  
**EĞİTİM BİLİMLERİ ANA BİLİM DALI**  
**EĞİTİMDE ÖLÇME VE DEĞERLENDİRME**  
**TEZLİ YÜKSEK LİSANS PROGRAMI**

**EĞİTİM ALANINDA MAKİNE ÖĞRENİMİ SINIFLANDIRMA**  
**ALGORİTMALARININ İNCELENMESİ**

**YÜKSEK LİSANS TEZİ**  
**CANAY CAN**

**Danışman:**  
**Doç. Dr. Hakan KOĞAR**

**Antalya, 2021**

## **DOĐRULUK BEYANI**

Yüksek lisans tezi olarak sunduĐum bu alıřmayı, bilimsel ahlak ve geleneklere aykırı dűşecek bir yol ve yardıma bařvurmaksızın yazdıĐımı, yararlandıĐım eserlerin kaynakalardan gösterilenlerden oluřtuĐunu ve bu eserleri her kullanımında alıntı yaparak yararlandıĐımı belirtir; bunu onurumla doĐrularım. Enstitü tarafından belli bir zamana baĐlı olmaksızın, tezimle ilgili yaptıĐım bu beyana aykırı bir durumun saptanması durumunda, ortaya ıkacak tüm ahlaki ve hukuki sonuçlara katlanacaĐımı bildiririm.

11/08/2021

**Canay CAN**

**AKDENİZ ÜNİVERSİTESİ**  
**EĞİTİM BİLİMLERİ ENSTİTÜSÜ MÜDÜRLÜĞÜNE**

**Canay CAN**'ın bu çalışması ..... tarihinde jürimiz tarafından **Eğitim Bilimleri** Ana Bilim Dalı **Eğitimde Ölçme ve Değerlendirme** Tezli Yüksek Lisans Programında **Yüksek Lisans Tezi** olarak **oy birliği** ile kabul edilmiştir.

**Başkan:** Doç. Dr. Güçlü ŞEKERCİOĞLU

(Akdeniz Üniversitesi/Eğitim Fakültesi/Eğitim Bilimleri Bölümü)

**Üye (Danışman):** Doç. Dr. Hakan KOĞAR

(Akdeniz Üniversitesi/ Eğitim Fakültesi/Eğitim Bilimleri Bölümü)

**Üye:** Dr. Öğr. Üy. Gizem UYUMAZ

(Giresun Üniversitesi/ Eğitim Fakültesi/Eğitim Bilimleri Bölümü)

**Yüksek Lisans Tezinin Adı:** Eğitim Alanında Makine Öğrenimi Sınıflandırma Algoritmalarının İncelenmesi

**ONAY:** Bu tez, Enstitü Yönetim Kurulunca belirlenen yukarıdaki jüri üyeleri tarafından uygun görülmüş ve Enstitü Yönetim Kurulunun ..... tarihli ve ..... sayılı kararıyla kabul edilmiştir.

## TEŞEKKÜR

Yüksek lisans tez süresince, bana her konuda ve her daim, ilgisi ve anlayışıyla destek olan değerli tez danışmanım **Doç. Dr. Hakan KOĞAR**'a,

Tez sürecinde değerli yorumlarını ve yardımlarını esirgemeyen **Doç. Dr. Bilal Barış ALKAN**'a,

Tez geliştirme ve değiştirme sürecinde değerli fikirleri, yorumları ve yönlendirmeleri ile destek olan **Dr. Öğr. Üy. Gizem UYUMAZ** ve **Doç. Dr. Güçlü ŞEKERCİOĞLU**'na,

Yüksek lisans eğitimi boyunca değerli bilgilerinden faydalanma imkânı bulduğum **Prof. Dr. Bayram BIÇAK** ve **Doç. Dr. Alper SİNAN**'a,

Hayatımın her anında yanımda olarak ve bana inanarak desteklerini biran olsun esirgemeyen canım annem **Nuray CAN**, canım babam **Mehmet CAN** ve yanlarındayken tamamlandığımı hissettiğim ablalarım **Zeynep YILMAZ** ve **Zülal CAN** ile canım kardeşim **Mustafa Nail CAN**'a,

Doğumuyla birlikte hayatıma neşe kaynağı ve umut olan biricik yeğenim **Leyla YILMAZ**'a,

Son olarak, adım ile anılacak bu çalışmada benimle birlikte anılmasını istediğim sevgili oğlum, canım dostum, kedim **Badem**'e

Sonsuz teşekkürlerimi sunarım.

## ÖZET

### EĞİTİM ALANINDA MAKİNE ÖĞRENİMİ SINIFLANDIRMA ALGORİTMALARININ İNCELENMESİ

Can, Canay

Yüksek Lisans, Eğitim Bilimleri Anabilim Dalı

Tez Danışmanı: Doç. Dr. Hakan KOĞAR

Ağustos 2021, 126 sayfa

Bu araştırma ile eğitim alanında büyük veri çalışmalarına temel oluşturmak için makine öğrenmesi algoritmalarından hangilerinin alanda kullanılabileceği tespit edilmeye çalışılmıştır. Bu doğrultuda çalışmada, Türkiye Öğrenci Değerlendirmesi gerçek veri seti ile öğretim elemanı kalitesinin araştırılması için öğretim elemanlarının performanslarının belirlenmesi ile ilişkisi olduğu düşünülen derse özel 28 soru ve 5 özellikten oluşan faktörler; Karar Ağacı, Rastgele Orman ve Yapay Sinir Ağları algoritmaları çerçevesinde incelenmiştir. Bahsi geçen bu üç algoritmanın sınıflandırma performansları doğruluk, duyarlılık, özgüllük, kesinlik ve F ölçütü ile araştırılmıştır. Çalışmanın verileri Kaliforniya Üniversitesi Makine Öğrenmesi Veri Havuzu'ndan hazır olarak alınmıştır. Veri seti, üç farklı öğretim elemanının 5820 Gazi Üniversitesi öğrencisi tarafından likert tipi ölçek ile değerlendirilmesinden oluşmaktadır. Verilerin analizi için R yazılımı ve R Studio ortamı kullanılmıştır. Araştırma sonucunda doğru sınıflama oranı %98.57 olan Karar Ağaçları algoritması, Rastgele Orman ve Yapay Sinir Ağlarına göre daha başarılı bulunmuştur. %98.03 ile karar ağacına çok yakın değere sahip olan Rastgele Orman algoritması ikinci en iyi performans gösteren algoritma olarak belirlenmiştir. %81.55 ile Yapay Sinir Ağları kabul edilebilir düzeyde performans göstermiş olsa da diğer algoritmalar ile karşılaştırıldığında düşük performanslı olarak kabul edilmektedir. Algoritmalar özgüllük oranı, duyarlılık oranı, kesinlik oranı ve F ölçütü çerçevesinde incelendiğinde yine Karar Ağaçları algoritması, Rastgele Orman ve Yapay Sinir Ağlarına göre daha başarılı bulunmuştur. Araştırmanın alt problemi olarak üç algoritmanın en önemli yordayıcısı ve manidarlık düzeylerinin karşılaştırılması araştırılmıştır. Sonuç olarak, üç algoritma için en önemli yordayıcı *sınıf* değişkeni olarak bulunmuş ve manidarlık düzeyi açısından üç algoritma arasında herhangi bir farklılığa rastlanılmamıştır.

**Anahtar Kelimeler:** Eğitim, C5.0 Karar Ağacı, Rastgele Orman, Yapay Sinir Ağları, R programlama, R Studio.

## **ABSTRACT**

### **INVESTIGATION OF CLASSIFICATION ALGORITHMS OF MACHINE LEARNING IN THE FIELD OF EDUCATION**

Can, Canay

Master's of Arts, Department of Education Sciences

Supervisor: Doç. Dr. Hakan KOĞAR

August 2021, 126 page

With this study, we tried to determine which of the machine learning algorithms can be used in the field to form the basis of big data studies in education. Accordingly, in the study, factors consisting of 28 questions and five features, which are thought to be related to the determination of the performance of instructors in order to investigate the real data set and the quality of instructor, the Student Assessment of Turkey, are: The Decision Tree was examined under the algorithms of Random Forest and Artificial Neural Networks. The classification performances of these three algorithms were researched with accuracy, precision, specificity, precision, and F criteria. The study's data are gathered from the University of California's Machine Learning Data Repository. The data set consists of evaluation of three different lecturers with an likert-type scale by 5820 Gazi University students. R software and R Studio media were used to analyze the data. As a result of the study, the decision trees algorithm was found to be more successful than the Random Forest and Artificial Neural Networks, with an accurate classification ratio of 98.57%. The Random Forest algorithm was identified as the second best performing algorithm with a value very close to the decision tree at 98.03%. The Artificial Neural Networks performed at an acceptable level (81.55%) but are considered low-performance compared to other algorithms. Again, the Decision Tree algorithm was found to be more successful compared to Random Forest and Artificial Neural Networks when analyzed by using criteria like specificity, sensitivity ratio, accuracy ratio and F. The study investigated a comparison of the most significant predictor and materiality levels of three algorithms as the sub-problem of the study. In conclusion, the most important predictive class variable for the three algorithms was found, and no difference in level of magnetism was encountered among the three algorithms.ass variable and no difference was ecountered in terms of the level of magnetism.



**Keywords:** Training, C5.0 Decision Tree, Random Forest, Artificial Neural Networks, R programming, R Studio.

## İÇİNDEKİLER

TEŞEKKÜR .....	iv
ÖZET .....	v
ABSTRACT .....	vii
ŞEKİLLER LİSTESİ.....	xi
TABLolar LİSTESİ .....	xii
KISALTMALAR LİSTESİ.....	xiii

### BÖLÜM I

#### GİRİŞ

1.1. Problem Durumu .....	1
1.2. Araştırmanın Amacı ve Problemleri .....	27
1.3. Araştırmanın Önemi .....	28
1.4. Araştırmanın Varsayımları .....	28
1.5. Araştırmanın Sınırlılıkları .....	28
1.6. Tanımlar .....	29

### BÖLÜM II

#### KAVRAMSAL ÇERÇEVE VE İLGİLİ ARAŞTIRMALAR

2.1. Çalışmada Kullanılan Algoritmalar .....	30
2.1.1. Karar Ağacı (Decision Tree).....	30
2.1.2. Rastgele Orman (Random Forest).....	35
2.1.3. Yapay Sinir Ağları (Artificial Neural Networks) .....	37
2.2. İlgili Araştırmalar .....	49
2.2.1. Uluslararası Araştırmalar .....	49
2.2.2. Ulusal Araştırmalar .....	58

### BÖLÜM III

#### YÖNTEM

3.1. Araştırmanın Modeli .....	61
3.2. Çalışma Grubunun Verileri .....	62
3.3. Veri Analizi .....	64

3.2.1.	Karar Ağaçları Uygulaması.....	66
3.2.2.	Rastgele Orman Uygulaması.....	66
3.2.3.	Yapay Sinir Ağları Uygulaması .....	67

## **BÖLÜM IV**

### **BULGULAR**

4.1.	Algoritmaların Performanslarının İncelenmesi.....	68
4.1.1.	Karar Ağacı .....	68
4.1.2.	Rastgele Orman .....	71
4.1.3.	Yapay Sinir Ağları .....	73
4.2.	Algoritmaların Sınıflandırma Performanslarının Karşılaştırılması.....	76
4.3.	Algoritmaların En Önemli Yordayıcılarının Belirlenmesi ve Manidarlıklarının Karşılaştırılması .....	78

## **BÖLÜM V**

### **SONUÇ, TARTIŞMA VE ÖNERİLER**

5.1.	Sonuç ve Tartışma.....	82
5.2.	Öneriler.....	86
5.2.1.	Uygulayıcılar İçin Öneriler .....	86
5.2.2.	Araştırmacılar İçin Öneriler .....	86
KAYNAKÇA .....		87
EKLER .....		101
Ek – 1. Veri Ön İşleme İçin Yazılan Kod.....		101
Ek- 2. C5.0 Karar Ağacı İçin Yazılan Kod.....		103
Ek – 3. Rastgele Orman İçin Yazılan Kod .....		105
Ek – 4. Yapay Sinir Ağları İçin Yazılan Kod.....		107
ÖZGEÇMİŞ .....		109
BİLDİRİM.....		111
İNTİHAL RAPORU .....		112

## ŞEKİLLER LİSTESİ

Şekil 1. 1. Veri Bilimi Venn Şeması .....	2
Şekil 1. 2. Büyük Veri Boyutları.....	3
Şekil 1. 3. Yapay Zekâ, Makine Öğrenmesi ve Derin Öğrenme Arasındaki İlişki. ....	4
Şekil 1. 4. Makine Öğrenmesi Yöntemlerinden Denetimli ve Denetimsiz Öğrenme .....	7
Şekil 1. 5. Denetimsiz Öğrenme.....	8
Şekil 1. 6. Denetimsiz Öğrenme Yöntemlerinden Kümeleme Tekniği.....	9
Şekil 1. 7. Kümeleme Grafiği.....	9
Şekil 1. 8. Denetimli Öğrenme.....	10
Şekil 1. 9. Basit Doğrusal Regresyon Modeli .....	12
Şekil 1. 10. İkili Sınıflandırma .....	14
Şekil 1. 11. Çoklu Sınıflandırma.....	14
Şekil 1. 12. İkili Sınıflandırma Süreci (a) ve Çok Sınıflı Sınıflandırma Süreci (b).....	16
Şekil 1. 13. <i>Eğitim ve Test Verisi</i> .....	18
Şekil 2. 1. Karar Ağacı .....	31
Şekil 2. 2. Karar Ağacı Kök, Dal ve Yaprak Dağılımı.....	31
Şekil 2. 3. Rastgele Orman.....	36
Şekil 2. 4. Biyolojik Sinir Ağı (a) ve Yapay Sinir Ağı (b).....	37
Şekil 2. 5. Tek Katmanlı Tek Nöronlu Yapay Sinir Ağı.....	38
Şekil 2. 6. Sigmoid Fonksiyonu .....	40
Şekil 2. 7. Aktivasyon Fonksiyonları .....	41
Şekil 2. 8. Tek Katmanlı Tek Nöronlu YSA Örnek Soru .....	42
Şekil 2. 9. Tek Katmanlı Çok Nöronlu Yapay Sinir Ağları .....	43
Şekil 2. 10. YSA Matematiksel Adımları .....	43
Şekil 2. 11. Tek Katmanlı İleri Beslemeli Ağlar.....	44
Şekil 2. 12. Çok Katmanlı İleri Beslemeli Ağlar .....	45
Şekil 2. 13. Geri Beslemeli Yapay Sinir Ağları .....	45
Şekil 2. 14. Yapay Sinir Ağlarında Eğitim.....	47
Şekil 3. 1. TÖD Veri Seti Korelasyonel İlişkisi .....	64
Şekil 4. 1. C5.0 Karar Ağacı Modeli.....	68
Şekil 4. 2. Yapay Sinir Ağları Modeli.....	73

## TABLULAR LİSTESİ

Tablo 1. 1. Karışıklık Matrisi .....	19
Tablo 1. 2. Kappa Değeri ve Yorumu .....	21
Tablo 3. 1. Türkiye Öğrenci Değerlendirmesi Veri Seti .....	62
Tablo 3. 2. TÖD Korelasyonel İlişki Sonuçları .....	65
Tablo 4. 1. Karar Ağacı Karar Kuralları .....	69
Tablo 4. 2. Eğitim Seti Hata Tablosu .....	70
Tablo 4. 3. Test Seti Hata Tablosu .....	70
Tablo 4. 4. OBB Hata Tablosu .....	71
Tablo 4. 5. Mtry Düzenlenmiş OBB Hata Tablosu.....	71
Tablo 4. 6. Eğitim Seti Hata Tablosu .....	72
Tablo 4. 7. Test Seti Hata Tablosu .....	72
Tablo 4. 8. Yapay Sinir Ağı Model Detayları .....	74
Tablo 4. 9. Test Seti Hata Tablosu .....	75
Tablo 4. 10. Doğruluk Tablosu .....	76
Tablo 4. 11. Özgüllük Tablosu.....	76
Tablo 4. 12. Duyarlılık Tablosu .....	77
Tablo 4. 13. Kesinlik Tablosu .....	77
Tablo 4. 14. F-Ölçütü Tablosu .....	77
Tablo 4. 15. Karar Ağacı İçin Değişkenlerin Önem Tablosu.....	78
Tablo 4. 16. <i>Rastgele Orman İçin Değişkenlerin Önem Tablosu</i> .....	79
Tablo 4. 17. Yapay Sinir Ağları İçin Değişkenlerin Önem Tablosu.....	80
Tablo 4. 18. Algoritmaların Manidarlık Düzeyinin Karşılaştırılması .....	81

## KISALTMALAR LİSTESİ

**$\Sigma$**  : Toplam Fonksiyonu

**$f(-)$**  : Aktivasyon Fonksiyonu

**AI**: Yapay Zekâ

**AID**: Automatic Detector (Otomatik Dedektör)

**ANN/YSA**: Yapay Sinir Ağları

**CART**: Classification and Regression Tree (Sınıflama ve Regresyon Ağacı)

**CHAID**: Chi-Square Automatic Interaction Detection (Ki-Kare Otomatik Etkileşim Algılama)

**DN**: Doğru Negatif

**DP**: Doğru Pozitif

**DT/KA**: Karar Ağaçları

**KGP**: Karanlık Gökyüzü Parkı

**KNN**: K-En Yakın Komşu

**ML/MÖ**: Makine Öğrenimi

**OEST**: Quick, Unbiased, Efficient Statistical Tree (Hızlı, Yansız, Etkili, İstatistiksel Ağaç)

**Pe**: Tahmin Değeri

**Po**: Gerçek Değeri

**RF/RO**: Rastgele Orman

**SVM**: Destek Vektör Makineleri

**TÖD**: Türkiye Öğrenci Değerlendirmesi

**W**: Ağırlık

**YN**: Yanlış Negatif

**YP**: Yanlış Pozitif

# BÖLÜM I

## GİRİŞ

Bu bölümde araştırmanın problem durumu, amacı, alt problemleri, önemi, varsayımları, sınırlılıkları ve tanımlar ele alınmıştır.

### 1.1.Problem Durumu

#### 1.1.1. Veri Bilimi (Data Science)

Veri, “herhangi bir işleme tabi tutulmadan, gözlem veya ölçüm yöntemleri ile ortamdan elde edilen her türlü değerdir” (Şeker, 2013, s. 22) ya da “tek başına anlam ifade etmeyen veya kullanılamayan, bununla birlikte enformasyona ve bilgiye temel oluşturan ilişkilendirilmeye, gruptandırılmaya, yorumlanmaya, anlamlandırılmaya ve analiz edilmeye gereksinim duyulan ham bilgi” (Yılmaz, 2009, s. 98) şeklinde tanımlanmaktadır.

Dhar (2013) “bilim” terimini, sistematik çalışma yoluyla elde edilen bilgi şeklinde tanımlamaktadır. Bilim, bilgiyi test edilebilir açıklamalar ve tahminler şeklinde oluşturan ve organize eden sistematik girişimdir. Bu nedenle veri bilimi, çıkarıma olan güvenimiz de dâhil olmak üzere, verileri ve dolayısıyla istatistikleri ya da organizasyonun sistematik çalışmasını, verilerin özelliklerini, analizini ve çıkarımdaki rolünü içeren bir odak noktasıdır. Veri bilimi, gerçek problemlerin verilerle isimlendirilmesi ve problemler için çözüm üretecek veri uygulamalarını geliştiren bilgisayar, istatistik, matematik ve bilişim bilimlerini içinde barındıran disiplinler arası çalışan bir alandır (Hamilton, 2015). Başka bir ifade ile büyük verinin karmaşıklığını çözümlmek için istatistik, makine öğrenimi, veri madenciliği ve veri tabanı teknolojilerini birleştirerek verileri analiz eden yöntemleri ifade etmektedir (Cavique, 2014). Şekilde 1.1’de yer alan venn şemasında Veri Bilimi’nin kullanıldığı alanlara yer verilmiştir.



Şekil 1. 1. *Veri Bilimi Venn Şeması*

**Kaynak:** (Şahin, 2019)

Veri Bilimi belirli ilkeler doğrultusunda veriden bilgi çıkarılmasını destekleyen ve yönlendiren bir dizi temel ilkedен oluşmaktadır. Analiz yöntemi ile ilişkisi en yakın olan kavram, veri madenciliğidir. Veri madenciliği, yüzlerce farklı algoritmayı kullanarak, bahsedilen ilkeler doğrultusunda ilgili teknolojiler ile veriden fiili olarak bilgi çıkarmaktadır. Bu ilke ve teknikler, birçok işlevsel alanda geniş bir şekilde uygulanmaktadır (Provost ve Fawcett, 2013).

Gerçek veri ile bir soruna çözüm olabilecek bilgileri elde etmek ve izlenmesi gereken süreci belirlemek için verilerin toplanması, hazırlanması, analizi ve analizlerin sonuçlarının yorumlanması gibi süreçleri ele alan veri bilimi, süreçlerin sürdürülebilmesi için büyük veri, veri madenciliği, derin öğrenme ve makine öğrenmesi teknoloji ve yöntemlerini kullanmaktadır (Hamilton, 2015).

### 1.1.2. Büyük Veri (Big Data)

Büyük veri, veri çağının başlangıcından itibaren var olan bir terimdir. Günümüzdeki anlamıyla ilk kez 1990'ların sonlarında kullanılmıştır. Diebolt (2000) tarafından, *Makroekonomik Ölçüm ve Tahmin için Büyük Veri Dinamik Faktör Modelleri* (Big Data Dynamic Factor Models for Macroeconomic Measurement and Forecasting) isimli bildiri bahsi geçen konu ile ilgili ilk akademik makaledir. Ancak büyük veri terimini ilk kullanan kişi Silicon Graphics (SGI)'nin baş veri bilimcisi John Mashey'dir. Mashey, 1990'ların sonlarında yaptığı



bir konuşmasında büyük verinin gelgit dalgasından söz etmiştir. Büyük veri çağı, hayallerin çok ötesinde, hızla genişleyen veri hacmi ile tanımlanan bir çağdır (Dean, 2014).

Direkt olarak kaynağından toplanan, birbirine benzemeyen verilerin analizi, işlenmesi ve depolanması ile ilgilenen büyük veri (Erl, Khattak ve Buhler, 2016); önceki süreçte ölçülmesi, saklanması, analiz edilmesi ve paylaşılması mümkün olmayan verilerin büyük çoğunluğunu, yeni sürecin başlaması ile birlikte ilgili veri ambarlarında saklamaya başlamıştır (Mayer-Schönberger ve Cukier, 2013).

Monino ve Sedkaoui (2016) büyük veri terimi, “işlem sırasında kullanılan veri hacminin kritik seviyeye ulaşması durumunda, yeni teknolojik depolama, işlem ve yöntemlerin kullanılmasıdır” şeklinde tanımlamaktadır. Fakat kendi içinde veri hızı, veri değeri, veri doğruluğu, veri çeşitliliği ve veri hızı (Şekil 1.2) ile bir sistem oluşturarak çalışan büyük veri için yalnızca fiziksel anlamda bir büyüklük olarak ele alınmamaktadır. Büyük veri aynı zamanda veri kümelerini açıklaması (Cackett, 2013), veri analizi sürecindeki önem ve etkisi açısından da “büyük” olarak ifade edilmektedir (Monino ve Sedkaoui, 2016).



Şekil 1. 2. *Büyük Veri Boyutları*

**Kaynak:** (Bayrakçı ve Albayrak, 2019)

### 1.1.3. Yapay Zekâ (Artificial Intelligence, AI)

Yapay zekâ, en basit ifade ile insan yapımı anlamına gelmektedir. Yaratılıştaki biçim, fiziksel konum veya mimari yapı yapay düzende karmaşıktır. Temelde silikon çiplerin

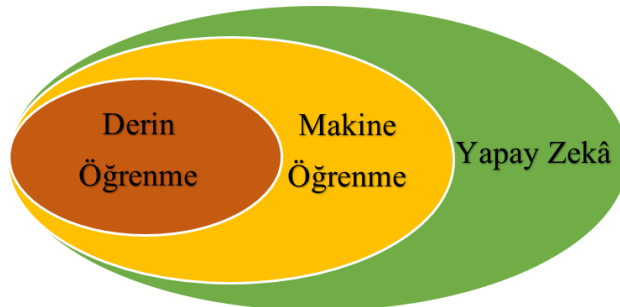
kullanıldığı cihazlarla çalışan bilgisayar programıdır. Günümüzde akıllı sistemlerde sıklıkla kullanılmaktadır (Wodecki, 2020).

Yapay zekâ “özerk” çalışabilmekte, uzaktan yönetilen akıllı sistemlerden yararlanabilmekte ve farklı nesnelere ile ağ kurarak etkileşime girebilmektedir. Kuantum bilgisayarlar tarafından gelişmiş işlemler gerçekleştirilebilmektedir. Veriler tıpkı genetik kod taşıyan DNA gibi kodlarla saklanabilmekte ve tıpkı protein sentezi gibi veri yapıları kullanılarak hesaplamalar yapılmaktadır (Wodecki, 2020).

Yapay zekâ alanının literatürde kabul edilmiş adıdır, ancak “yapay zekâ” terimi kafa karışıklığına neden olabilmektedir. Çünkü yapay zekâ birçok insan tarafından gerçek zekânın tersi olarak yorumlanmaktadır (Poole ve Mackworth, 2010). Yapay zekâ programı “Akıllı Ajan” olarak da adlandırılmaktadır (Das, Dey, Pal ve Roy, 2015). Bu nedenle, akıllı hareket eden hesaplama ajanlarının sentezini ve analizini inceleyen bir alan olarak tanımlanmaktadır (Poole ve Mackworth, 2010). Ajanlar, algılayıcıları (sensörleri) aracılığıyla bir ortamın durumunu belirleyebilmekte ve ardından bir mekanizmayı veya sistemi kontrol eden veya hareket ettiren sistemler (aktüatörleri) ile durumu etkileyebilmektedir. Yapay zekânın önemli yönü, algılayıcılardan elde edilen girdilerin kontrol sistemine nasıl çevrildiğini, başka bir ifade ile algılayıcıların kontrol sistemi ile nasıl eşlendiğini ima eden bir ajan kontrol politikasına sahip olmasıdır. Bu politika ise, bir fonksiyon ile gerçekleşmektedir (Das ve diğerleri, 2015).

Yapay zekânın nihai amacı, makinelerde insan zekâsına benzer bir zekâ geliştirmektir. Bunun için insan beyninin nasıl öğrendiğini taklit eden algoritmalar geliştirilmiştir. Makine öğrenmesi açık programlama olmadan makinelerin insan benzeri zekâ kazanmasını sağladığı için büyük önem taşımaktadır (Das ve diğerleri, 2015).

Şekil 1.3’te yapay zekânın kapsadığı alanlara yer verilmiştir. Şekilde derin öğrenmenin makine öğreniminin alt alanı olduğu, yapay zekânın ise bu iki alanı kapsadığı gözlenmektedir.



Şekil 1. 3. Yapay Zekâ, Makine Öğrenmesi ve Derin Öğrenme Arasındaki İlişki.

**Kaynak:** (Bingöl, Ormecioglu ve Arzu, 2020)

#### 1.1.4. Makine Öğrenmesi (Machine Learning)

Makine öğrenmesi yapay zekâ, istatistik, olasılık, felsefe, psikoloji ve sinirbilimi (nörobilim) ile çalışabilen çok disiplinli bir alandır. Ulusal ya da uluslararası düzeyde oluşabilecek sorunları çözmek için, istatistiği doğru ve faydalı bir şekilde kullanarak model oluşturmaktadır. Bilgisayar sistemlerinin insan zihnini taklit etmeyi öğrenmesinin gerekip gerekmediğini keşfetme isteği ve taklitlerin temel istatistiksel hesaplamalar ile incelenmesi üzerine birçok teori oluşturulmuş (Domingos, 2012) ve böylece geniş kapsamlı uygulamalar yapılarak, Makine Öğrenmesi bilgisayar biliminin en hızlı büyüyen alanlarından biri haline gelmiştir (Shalev-Shwartz ve Ben-David, 2014).

Makine öğrenmesi hesaplamalı zekânın diğer bir ifade ile yapay zekânın önemli bir alt alanıdır. Bilgisayara veri ya da deneysel gözlem ile elde edilmiş gözlemlerin öğretilmesi yeteneğidir (Peng, 2013). Genel olarak sistemin performansını arttırmak veya doğru tahminler yapmak için deneyimini kullanan bir hesaplama yöntemidir (Mohri, Rostamizadeh ve Talwalkar, 2018).

1946 yılında ilk ENIAC aygıtı geliştirilmiştir. Böylece insan düşüncesinin ve öğrenmesinin makinelere mantıksal olarak işlenebileceği sonucuna ulaşılmıştır. Makinelere öğrenmeyi öğretme çalışmaları 1950’de Alan Turing tarafından devam ettirilmiştir (Domingos, 2012). Turing (1950) çalışmalarında “makinelere öğrenebilir mi?” sorusuna cevap aramış ve Turing Testi’ni geliştirmiştir. Turing testi, makine ile iletişim kurulup kurulmadığı, ayırt edebilme yeteneğinin var olup olmadığı gibi insana ait öğrenme becerilerini makinenin kazanıp kazanmadığının incelenmesi fikrine dayanmaktadır. Makinelerin öğrenebilirliğinin test edildiği çalışmalara 1952 yılında Arthur Samuel (IBM) dâhil olmuştur. Uluslararası dama şampiyonunu mağlup etmesi için ilk spor-oyun yazılımını yazmıştır. 1957 yılında Psikolog Frank Rosenblatt, uygulama içindeki karmaşık sorunları çözmek için verileri birbirine bağlayan algılayıcı sistemi icat etmiştir. Bu icadı ile derin sinir ağlarının temelini atmıştır. 1967’de örüntü odaklı çalışmaların geliştirilmesiyle en yakın komşu algoritması tasarlanmıştır (Domingos, 2012). 1981’de Gerold De Jong çalışma alanına ilişkin önceki bilgileri denetimli öğrenme ile ele alarak çalışmalarını sürdürmüştür. 90’lı yılların başlarında ise, Bilgisayar Bilimi ve İstatistik’in kesişimi olan makine öğrenimi yeniden popülerlik kazanmıştır (Solanki ve Dhankar, 2017).

Makine öğrenimi, eldeki verilerden bilgi elde etmek için birtakım hesaplama metotları kullanmaktır. Robotik ve bilgisayar oyunları, el yazısı ve konuşma tanıma, doğal dil işleme gibi

oldukça fazla uygulama alanına sahiptir (Hsieh, 2009). Makine öğrenimi ile insanların üstesinden gelemedikleri ya da yetersiz kaldıkları problemler kolaylıkla çözülebilmektedir. Gerçek hayat sürekli bir değişim içerisindedir ve hayat ihtiyaçları doğrultusunda evrilmektedir. Birçok yazılım programına her yıl birtakım özellikler eklenerek ya da çıkartılarak tekrar piyasaya sürülmektedir. Bu bilgisayara indirilen yazılım programlarının değişime ya da güncellemeye ihtiyacı olduğunu göstermektedir. Makine öğrenmesi ise, onu içeren sistemler de kolaylıkla değişikliklere uyum sağlayabilmektedir (Shalev-Shwartz ve Ben-David, 2014). Bu gibi özelliklerinden dolayı makine öğrenimi son yıllarda bilgi teknolojisinin temel dayanağı haline gelmiştir (Smola ve Vishwanathan, 2008).

Makine öğrenmesi bilgisayarlardaki bir problemi çözmek için algoritmalara ihtiyaç duymaktadır. Algoritma, girdiyi çıktıya dönüştürmek için gerçekleştirilmesi gereken bir dizi talimattır. Örneğin, bir sıralama çalışması için algoritma geliştirilebilir. Bunun için girdi bir dizi sayıdır ve çıktı onların sıralı listesidir. Aynı görev için çeşitli algoritmalar kullanılabilir. Algoritmalarından hangisinin seçileceğine en verimli olanı araştırılarak karar verilmektedir. (Alpaydın, 2020). Makine öğrenmesi sürekli değişen ve gelişen bir alandır (Boucheron ve Tagliaferri, 2019).

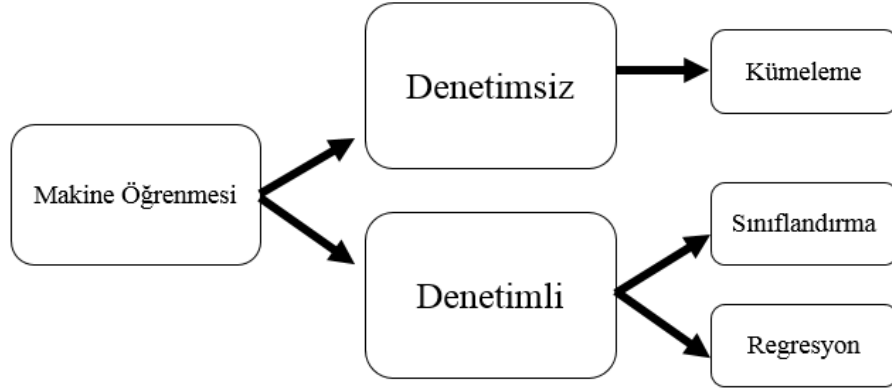
#### **1.1.4.1. Makine Öğrenmesi Yöntemleri**

Makine öğreniminde görevler genellikle geniş kategorilere ayrılmaktadır. Bu kategoriler, öğrenmenin nasıl alındığına veya geliştirilen sisteme öğrenmeyle ilgili geri bildirim nasıl verildiğine dayanmaktadır (Boucheron ve Tagliaferri, 2019). Sistemsel olarak farklı işleyiş gösteren birçok yöntem bulunmaktadır. Bunlar; Yarı Denetimli Öğrenme, Transdüktif Akıl Yürütme, Çevrimiçi Öğrenme, Pekiştirmeli Öğrenme ve Aktif Öğrenme olarak ayrılmaktadır (Mohri, Rostamizadeh ve Talwalkar, 2012). En yaygın kullanılan makine öğrenimi yöntemlerinden ikisi, insanlar tarafından etiketlenen örnek girdi ve çıktı verilerine dayalı algoritmaları eğiten denetimli öğrenme ve algoritmanın girdisi içinde yapı bulmasını sağlamak için etiketli veri içermeyen denetimsiz öğrenmedir (Boucheron ve Tagliaferri, 2019).

#### **Denetimli ve Denetimsiz Öğrenme**

Makine öğreniminde kullanılacak birçok farklı algoritma vardır. Veri setine göre hangi algoritmanın kullanılacağına karar verilmektedir. Bir veri setine birden fazla algoritma uygulanabilmektedir. Hangisinin daha iyi olduğunun belirlenmesi için sonuçlar karşılaştırılmaktadır. Makine öğrenmesi adından da anlaşılacağı üzere öğrenme üzerine kurulu bir sisteme sahiptir. Algoritmalar karakteristik özelliklerinden dolayı genel olarak iki öğrenme

türü üzerinden ilerlemektedir. Bunlar; Şekil 1.4'te görseli yer alan denetimli ve denetimsiz öğrenme yöntemleridir (Bell, 2020).



Şekil 1. 4. *Makine Öğrenmesi Yöntemlerinden Denetimli ve Denetimsiz Öğrenme*

**Kaynak:** (Ay, 2020)

#### i. **Denetimsiz veya Gözetimsiz Öğrenme (Unsupervised Learning)**

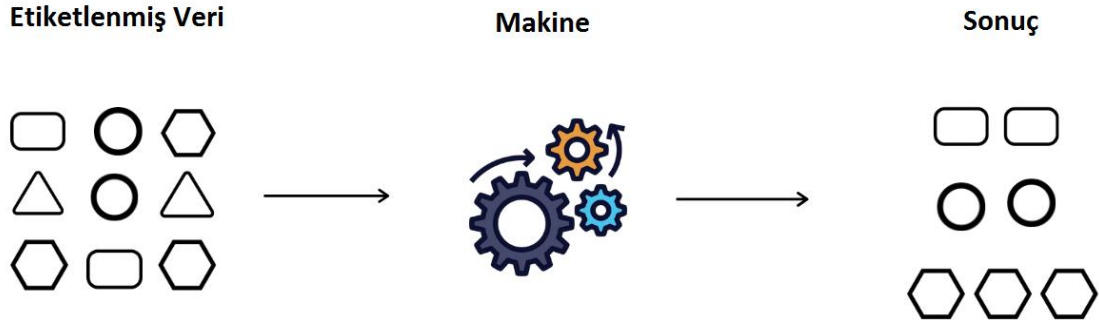
Denetimsiz öğrenme, etiketlenmemiş verilerin istatistikler hakkında denetimsiz bilgi edinme isteğidir. Amacı, istatistikleri keşfetmek ve ayrıntılar arasındaki benzerlikleri bulmaktır (Brownlee, 2015). Denetimsiz öğrenme bir kontrol sistemi tarafından kontrol edilmemektedir. Yöntem, verilerin belirli bölümünü ya da sıkıştırılmış halini bir yargıya varmak amacıyla kullanılabilir hale getirmektedir (Shalev-Shwartz ve Ben-David, 2014).

Denetimsiz öğrenmede sisteme sadece girdi değişkeni verilmektedir. Veri setinde herhangi bir çıktı ya da etiketleme bilgisi yer almamaktadır. Öğrenme sürecinde denetimden sorumlu bir denetleyici bulunmamaktadır. Sistem, değişkenler arasındaki ilişkiyi incelemekte ve duruma göre kümeleme gibi işlemler doğrultusunda çalışmaya başlamaktadır (Mohri, Rostamizadeh ve Talwalker, 2012). Denetimsiz öğrenme modeli için çözüme ulaşma yolculuğunda doğru ya da yanlış söz konusu değildir. İşlenen verinin örüntüsünü ortaya çıkarmak ve anlamlı bir bütünlük sağlamak amacı vardır (Bell, 2020). Veri setinde çıktı bilgisi yer almadığı için denetimsiz öğrenme için öğrenme performansını belirlemek oldukça zor bir şekilde gerçekleşmektedir (Mohri, Rostamizadeh ve Talwalker, 2012).

Örneğin, büyük bir turizm şirketi düzenlediği her organizasyonun ardından müşterilerine bir anket uygulamaktadır. Anket verileri ile denetimsiz öğrenme algoritmaları

kullanılarak müşteriler kümelenebilmektedir. Elde edilen sonuçlar ile müşterilere özel promosyon gibi sürprizler yapılarak finansal açıdan şirket kendini büyütmektedir (Bell, 2020).

Şekil 1.5'te denetimsiz öğrenme yöntemi görselleştirilmiştir.



Şekil 1. 5. Denetimsiz Öğrenme

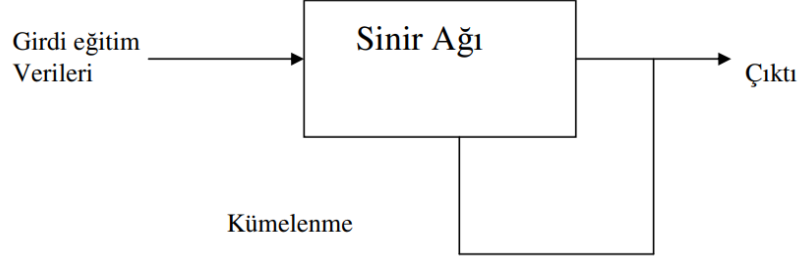
**Kaynak:** (Raj, 2021)

Şekil incelendiğinde, denetimsiz öğrenme yöntemi için makineye etiketlenmemiş verilerin yüklendiği, sonuç olarak ise verilerin benzerliklerine göre kümelendiği gözlenmektedir.

## **i.i. Kümeleme**

Kümeleme, benzer özellikleri paylaşan bir grup değişkenin organize edilmesi şeklinde tanımlanmaktadır. Denetimsiz öğrenme modellerinden birisidir ve bu yüzden makinenin önceden öğrenmesi gereken bir eğitim seti bulunmamaktadır. Denetimsiz öğrenmede amaç, belirli bir veri kümesi içindeki yapıyı bulmaktır. Kümeleme yöntemi geniş bir ağ oluşturmaktadır ve böylece aralarında seçim yapılabilecek çok sayıda algoritma oluşturmaktadır (Bell, 2020).

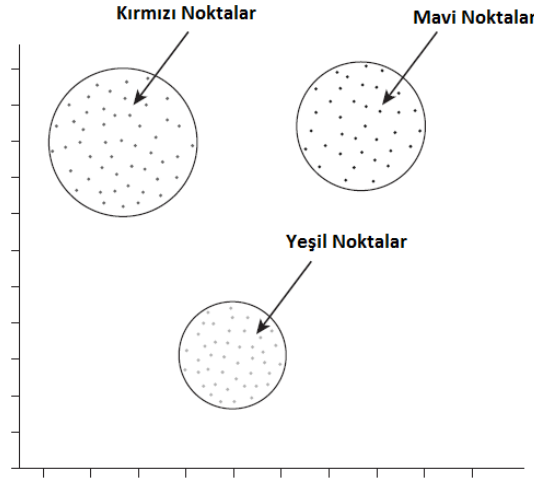
Örneğin, pazarlamacılar tarafından farklı pazarlama programları geliştirilir. Burada amaç müşteri kaçırmamaktır. Bu nedenle sektör, kümeleme yöntemi ile müşterileri ilgi alanlarına göre gruplara ayırarak çalışmalarını sürdürmektedir (Bell, 2020). Şekil 1.6’da kümeleme yönteminin görseli yer almaktadır.



Şekil 1. 6. Denetimsiz Öğrenme Yöntemlerinden Kümeleme Tekniği

**Kaynak:** (Karakuzu, 1998)

Şekil 1.7’de ise denetimsiz öğrenme yöntemlerinden kümeleme tekniği için bir çalışmanın grafiğine yer verilmiştir.



Şekil 1. 7. Kümeleme Grafiği

**Kaynak:** (Bell, 2020)

## ii. Denetimli veya Gözetimli Öğrenme (Supervised Learning)

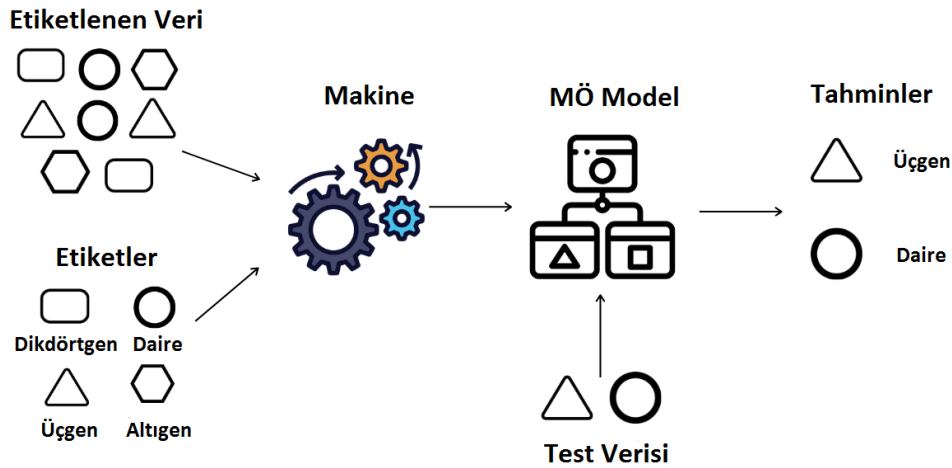
Denetimli öğrenmede bilgisayara etiketlenmiş örnek girdiler verilmektedir. Yöntemin amacı algoritmanın hatalarını bulmak için gerçek çıktı değerlerinin “öğretilen” çıktılarla karşılaştırılması ve makinenin öğrenmesini sağlayarak modelin buna göre değiştirilmesi esasına

dayanmaktadır. Bu nedenle denetimli öğrenme, etiketlenmemiş veriler üzerindeki etiket değerini tahmin etmek için kalıpları kullanmaktadır (Boucheron ve Tagliaferri, 2019).

Denetimli öğrenme yöntemi ile girdi ve çıktı verileri denetim ve gözetim çerçevesinde makineye aktarılmakta ve eldeki bilgilerden anlamlı sonuç elde edilmektedir. Bu öğrenme modelinin uygulanabilmesi için girdi ve çıktıları içeren bir eğitim veri setine ihtiyaç duyulmaktadır. Elde edilen eğitim veri seti denetmen aracılığıyla makineye öğretilmektedir. Böylece makine, veriler arasındaki ilişkiyi öğrenecek ve daha önce tanıtılmamış bir örneklem varsayımında bulunabilecektir. Burada denetiminin görevi, sisteme girdileri ve amaçlanan çıktıları öğretmektir. Denetimli öğrenmede, eğitim verisi ile sisteme deneyim kazandırılmakta, test verisi ile algoritmanın performansı değerlendirilmektedir (Shalev-Shwartz ve Ben-David, 2014).

Örneğin, Twitter üzerinde “tivit sınıflandırması” yapmak isteyen bir kişi öncelikle sisteme tivitleri girdi olarak göstermekte ve ardından amaçlanan çıktıları tanıtmaktadır. Amaçlanan çıktı, sistemin hangi konu altında sınıflandırma yapacağıdır. Makine sistemi girdiler ve çıktılar arasındaki gerekli ilişkilendirmeyi yaptıktan sonra öğrendiği eğitim seti doğrultusunda sınıflandırma yapmaktadır (Bell, 2020).

Şekil 1.8’de denetimli öğrenme yöntemi görselleştirilmiştir.



Şekil 1. 8. Denetimli Öğrenme

**Kaynak:** (Raj, 2021)



Şekil incelendiğinde, denetimli öğrenme yöntemi için veri seti, eğitim ve test verisi olarak ikiye ayrılmaktadır. Eğitim verisi ile model oluşturulduğu, test verisi ile modelin tahmin performansının kontrol edildiği gözlenmektedir.

Denetimli öğrenme, istatistiksel olarak olası gelecekteki olayları tahmin etmek için geçmiş verileri kullanmaktadır. Bu bilgi makine öğrenmesi yöntemleri için en yaygın kullanım olarak literatürde yer almaktadır (Boucheron ve Tagliaferri, 2019). Denetimli öğrenme algoritmaları, regresyon ve sınıflandırma yöntemlerinde kullanılmaktadır (Bontempi, Taieb ve Le Borgne, 2013).

### **ii.i. Regresyon**

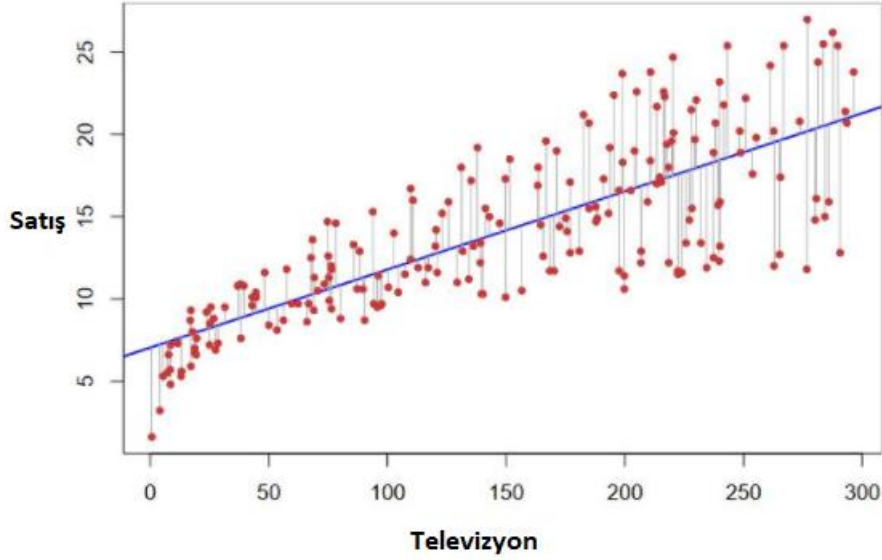
Regresyon, gözlemlenen  $x$  ölçümlerine dayalı olarak tek değişkenli  $w$  dünya durumunu tahmin etmek olarak tanımlanmaktadır (Prince, 2012). Bir başka ifade ile *bir değişkene ilişkin ölçümlerin grup ortalamasına doğru çekilmesi* şeklinde tanımlanmaktadır. Regresyon analizi, iki ya da daha fazla değişkenden birinin bağımlı diğerinin bağımsız olduğu durumlarda değişkenler arasındaki ilişkinin matematiksel eşitlik ile açıklanmasıdır (Büyüköztürk, Çokluk ve Köklü, 2019).

Regresyon analizinin amacı; bağımlı değişken ile bağımsız değişken ya da değişkenler arasındaki ilişkiyi açıklamak, bağımsız değişken ya da değişkenler ile bağımlı değişkeni tahmin etmek, bağımlı değişkende gözlenen değişimleri determinasyon katsayısı ile açıklamak ve birden fazla bağımsız değişken olduğu durumda değişkenlerin bağımlı değişken üzerindeki önemliliklerini belirlemektir. İki ya da daha fazla değişken arasındaki neden sonuç ilişkisinin tespit edilmeye çalışılması amacıyla kullanılan regresyon analizi, basit ve çok değişkenli regresyon olarak ikiye ayrılmaktadır. Çalışma bir bağımlı ve bir bağımsız değişkenden oluşuyorsa buna “doğrusal regresyon”, birden çok bağımsız değişken kullanılıyorsa buna “çok değişkenli regresyon” denilmektedir (Ünver ve Gamgam, 1986; Howell, 1987).

### **Doğrusal Regresyon (Linear Regression)**

Doğrusal regresyon, bir bağımsız ve bir bağımlı iki değişken arasındaki doğrusal ilişki şeklinde tanımlanmaktadır. Değişkenler arasında doğrusal ya da doğrusal olmayan ilişki olabilmektedir (Büyüköztürk, Çokluk ve Köklü, 2019).

Şekil 1.9’da bağımlı (satış) ve bağımsız (televizyon) değişkenler arasında doğrusal ve doğrusal olmayan ilişki türlerine yönelik görseller yer almaktadır (Efe, Bek ve Şahin, 2000).



Şekil 1. 9. *Basit Doğrusal Regresyon Modeli*

**Kaynak:** (James, Witten, Hastie ve Tibshirani, 2013)

Bağımlı (Y) ve bağımsız (X) değişkenler arasındaki doğrusal ilişkinin matematiksel eşitliği;

$$Y_i = \alpha + \beta X_i + e_i$$

şeklindedir (Efe, Bek ve Şahin,2000).

### Çok Değişkenli Regresyon

Çok değişkenli regresyon, iki veya daha fazla değişken ile bir bağımlı değişken arasındaki doğrusal ilişki şeklinde tanımlanmaktadır (Büyüköztürk, Çokluk ve Köklü, 2019).

Matematiksel formül;

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_2 + \dots + \beta_p X_{pi} + \epsilon_i$$

şeklindedir (Efe, Bek ve Şahin,2000).

## ii.ii. Sınıflandırma

Sınıflandırma, bilgisayar programına verilen verinin makine tarafından öğrenilmesi ve yeni gözlemlerin kolaylıkla sınıflandırılması için model tasarlayan denetimli öğrenme yaklaşımıdır (Bell, 2020). Sınıflandırma, tahmini bir sınıflandırma çıktısı veren bir bilgisayar programı olarak düşünüldüğünde matematiksel terimlerle, sınıflandırıcı örnekleri etiketlere eşleyen bir işlev olarak tanımlanabilmektedir (Schapire ve Freund, 2012).

Sınıflandırma adımları regresyondan farklı ilerlemektedir. Sınıflandırma için bir veri noktasına sürekli bir değer atamak yerine sınıf atanmaktadır. Bu nedenle, veri noktası başına daha az etiket seçeneği olduğundan görev daha basittir, ancak istenilen sınıfı elde etmek için her sınıf hakkındaki bilgileri bir şekilde hesaba katmak gerektiğinden süreç daha karmaşık ilerlemektedir (Nasiriany ve diğerleri, 2019).

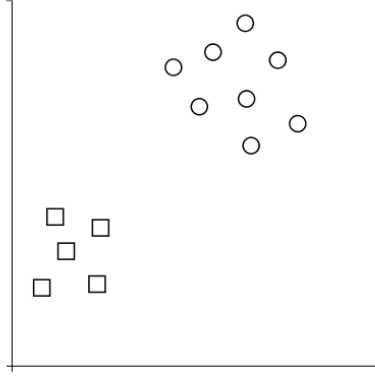
Sınıflandırma yapabilmek için veri, eğitim ve test seti olarak ikiye ayrılmaktadır. Makinenin eğitimi sırasında algoritma, eğitim örnekleri adı verilen etiketli örneklerden oluşan bir eğitim seti kullanmaktadır. Eğitim setinin çıktısı sınıflandırma kuralını vermektedir. Sınıflandırma çalışmalarında belirli bir sınıflandırıcının kalitesini değerlendirmek için hata oranı, yani yanlış sınıflandırma yapma sıklığını ölçülmektedir. Bunu yapmak için bir test setine ihtiyaç duyulmaktadır. Sınıflandırıcı, test örnekleri üzerinde bir değerlendirme gerçekleştirmektedir. Yanlış sınıflandırmaların yapıldığı örneklerden oluşan gruplara sınıflandırıcının test hatası denilmektedir. Benzer şekilde, eğitim kümesindeki hata oranına da eğitim hatası denilmektedir. Doğru tahmin oranına ise, doğruluk denilmektedir. Sınıflandırıcının performansı incelenirken test sonuçları dikkate alınmamaktadır. Öte yandan eğitim ve test seti arasındaki ilişki de göz ardı edilmemektedir (Schapire ve Freund, 2012).

Sınıflandırma, ikili sınıflandırma (bi-class) ya da çoklu sınıflandırma (multi-class) olarak iki başlıkta incelenmektedir (Jordan, Kleinberg ve Schölkopf, 2006).

### İkili Sınıflandırma

Bir lineer diskriminant fonksiyonu,  $w$ 'nin bir ağırlık vektörü ve  $w_0$ 'ın bir sapma değeri olduğu (istatistiksel anlamdaki sapma değeri ile karıştırılmamalıdır) girdi vektörünün lineer bir fonksiyonunu alarak elde edilmesinden oluşmaktadır. Sapma değerinin negatifine bazen eşik denilmektedir. Bir girdi vektörünün  $x$  olduğu durumda,  $y(x)=0$  ise  $C_1$  sınıfına, aksi takdirde  $C_2$  sınıfına atanmaktadır (Jordan ve diğerleri, 2006). Veri setinde iki etiketin olduğu sınıflandırma çalışmalarına ikili sınıflandırma denilmektedir (Bell, 2020).

Şekil 1.10’da ikili sınıflandırma yöntemi görselleştirilmiştir.



Şekil 1. 10. *İkili Sınıflandırma*

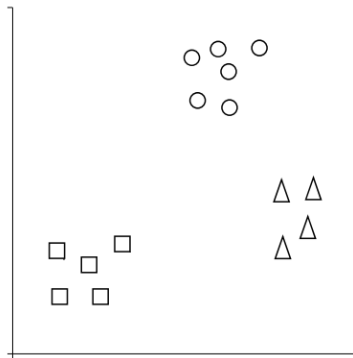
**Kaynak:** (Bell, 2020).

Şekil incelendiğinde veri setinin kare ve daire şeklinde sınıflandırıldığı gözlenmektedir. Burada dairelerin “evet”i, karelerinse “hayır”ı temsil ettiği varsayılmaktadır. Amaç ise, veri setinin hangilerinin evet (daire) sınıfına ait olduğu, hangilerinin de hayır (kare) sınıfına ait olduğunu belirlemektir. İşte bu tür çalışmalarda iki sınıfa ulaşılmaktadır. Dolayısıyla çalışma ikili sınıflandırma olarak adlandırılmaktadır (Bell, 2020).

### Çoklu Sınıflandırma

Çok sınıflı sınıflandırma, ikili sınıflandırmanın doğal bir uzantısıdır. Burada amaç, örneklerin her birine ayrı bir etiket atamaktır. İkili sınıflandırma çalışmalarında ele alınan ikili sınıflandırıcı ile çok sınıflı sınıflandırma problemi çözülebilmektedir (Daume, 2017).

Şekil 1.11’de çoklu sınıflandırma yöntemi görselleştirilmiştir.



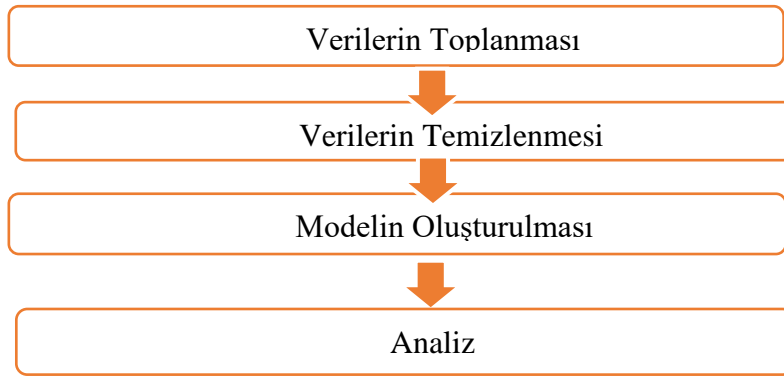
Şekil 1. 11. *Çoklu Sınıflandırma*

**Kaynak:** (Bell, 2020)

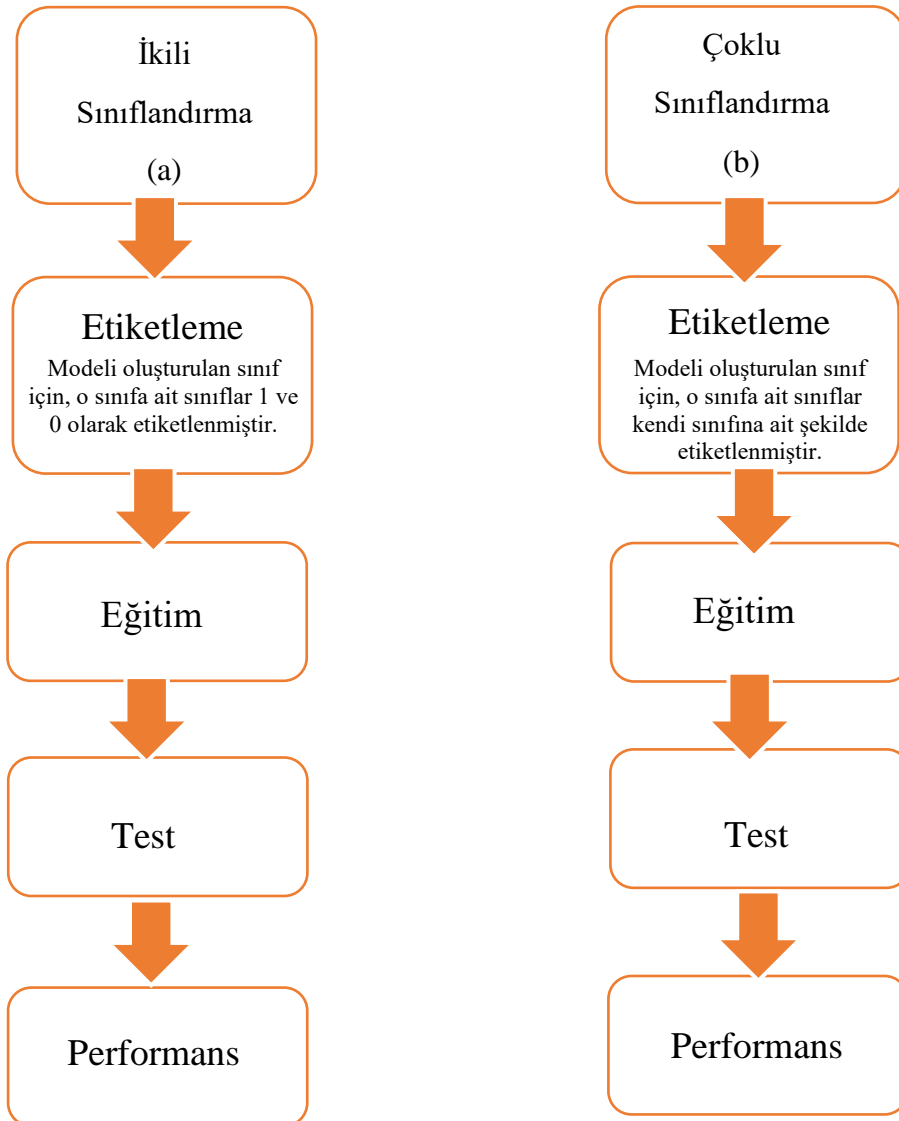
Şekil 1.11’de olduğu gibi veri setinde üç adet (veya daha fazla) sınıf var ise bu çok sınıflı bir sınıflandırma problemi ile karşılaşıldığını göstermektedir ve bu tür sınıflandırma çalışmaları çoklu sınıflandırma olarak adlandırılmaktadır (Bell, 2020).

Sınıflandırma işleminin gerçekleştirilebilmesi için birtakım süreçlerin izlenmesi gerekmektedir. Şekil 1.12’de genel olarak ve özel olarak ikili ve çoklu sınıflandırma sürecinin basamakları yer almaktadır.

## Sınıflandırma Süreci



## Örnek İkili ve Çoklu Sınıflandırma Süreci



Şekil 1. 12. İkili Sınıflandırma Süreci (a) ve Çok Sınıflı Sınıflandırma Süreci (b)

**Kaynak:** (Şahin ve Chouseinoglou, 2019)

## **Algoritma Performansını Değerlendirme Yöntemi**

Genellikle bir algoritmanın çıktı tahmin edicisinin gerçek değere yakın sonuçlar vermesi, yani bir algoritmadan daha iyi tahmin başarısı göstermesi beklenmektedir. Doğruluğu yüksek, gerçek bir tahmin için algoritmanın çıktı başarısı doğrulama seti adı altında eğitim verilerinin bir kısmı kullanılarak elde edilmektedir. Hold out yöntemi, algoritma performansının değerlendirilmesi için sıklıkla kullanılan yöntemlerden birisi olarak literatürde yer almaktadır (Shalev-Shwartz ve Ben-David, 2014).

### **Hold Out Yöntemi**

Hold out yöntemi, mevcut veri kümesini rastgele iki bölüme ayırmaktan oluşmaktadır. Bölümlerden biri modeli elde etmek için, diğer ise onu test etmek için kullanılmaktadır (Torgo, 2011). Bir başka ifade ile yöntem, bir tahmin edicisinin gerçek hatasını tahmin etmenin en basit yolu şeklinde tanımlanmaktadır. Burada eğitim setinden bağımsız olarak ek bir örnek veri seti belirlenmekte ve buna doğrulama seti denilmektedir. Doğrulama seti, deneysel hata tahmin edici olarak kullanılmaktadır. Bir eğitim setini örneklemek ve ardından bağımsız bir doğrulama setini örneklemek, rastgele örnek setini rastgele iki parçaya bölmekle eşdeğerlik taşımaktadır. Bölümlerden biri eğitim için diğeri ise doğrulama için kullanılmaktadır. Bu nedenle, doğrulama kümesine genellikle bekleme kümesi denilmektedir (Shalev-Shwartz ve Ben-David, 2014).

Hold out yöntemi yaygın olarak, veri setinin %70'i eğitim kümesi ve %30'u test kümesi olacak şekilde kullanılmaktadır. Veri setinin küçük olduğu durumlarda test kümesinin çok küçük olacak olması ya da eğitim setinden çok fazla veri çıkma olasılığından dolayı bahsi geçen yöntem genellikle çok büyük veri setleri için kullanılmaktadır (Torgo, 2011).

### **Eğitim Veri Seti ve Test Veri Seti**

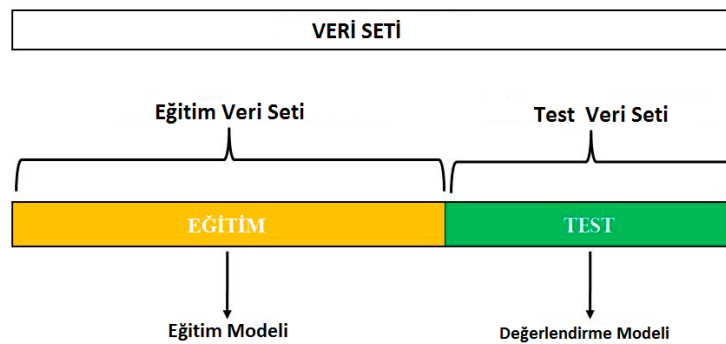
Sınıflandırma problemleri için, bir sınıflandırıcının performansı hata oranı açısından ölçülmektedir. Sınıflandırıcı, her örneğin sınıfını tahmin eder: Eğer doğruysa, bu bir başarı olarak sayılmaktadır; doğru değilse, bu bir hata olarak kabul edilmektedir. Hata oranı, bir dizi örnekte yapılan hataların oranıdır ve sınıflandırıcının genel performansını ölçmektedir (Witten, Friedman ve Simon, 2011).

Eğitim setindeki hata oranının gelecekteki performansın iyi bir göstergesi olması muhtemel değildir. Çünkü eğitim seti sınıflandırıcısı aynı eğitim verilerinden oluşmaktadır ve bu durumda, verilere dayalı herhangi bir performans tahmininin sonucu muhakkak iyimser olacaktır. Eğitim verileri üzerindeki hata oranı, yeniden ikame hatası olarak adlandırılmaktadır.

Çünkü eğitim örnekleri, onlardan oluşturulan bir sınıflandırıcıda yeniden ikame edilerek hesaplanmaktadır. Yeni verilerdeki gerçek hata oranının güvenilir bir tahmincisi olmamasına rağmen, yine de bilmek genellikle yararlı olmaktadır. Bir sınıflandırıcının yeni veriler üzerindeki performansını tahmin etmek için, sınıflandırıcının oluşumunda hiçbir rol oynamayan bir veri kümesindeki hata oranını değerlendirmesi gerekmektedir. Bu bağımsız veri kümesine test kümesi denilmektedir (Witten ve diğerleri, 2011).

Veri kümeleri eğitimden önce birleştirilirse test verileri üzerindeki performans muhtemelen tamamen farklı bir durumda olacak ve gelecekteki veriler üzerindeki performansın iyi bir göstergesi olmayacaktır. Test verilerinin herhangi bir şekilde model oluşturmak için kullanılmaması önem arz etmektedir. Bu verilerin hiçbiri tek başına gelecekteki hata oranının bir tahminini belirlemek için kullanılamaz. Bu gibi durumlarda genellikle üç veri kümesinden bahsedilmektedir: Eğitim verileri, doğrulama verileri ve test verileri (Witten ve diğerleri, 2011).

Eğitim verileri sınıflandırıcı oluşturmak için bir veya daha fazla öğrenme şeması tarafından kullanılmakta, doğrulama verileri ise bu sınıflandırıcının parametrelerini optimize etmek veya belirli bir tanesini seçmek için kullanılmaktadır. Test verileri ise nihai, optimize edilmiş yöntemin hata oranını hesaplamak için kullanılmaktadır. Test örneği ne kadar büyük olursa, hata tahmini o kadar doğru olur. Bu eğitim, doğrulama ve test için kullanılabilir veri miktarını sınırlamakta ve sınırlı bir veri kümesinden en iyi şekilde nasıl yararlanılacağı göstermektedir. Veri kümesinin belirli bir miktarı test için, geri kalanı ise eğitim için kullanılmaktadır (Witten ve diğerleri,2011; Molnar, 2020). Şekil 1.13'te veri setinin eğitim ve test seti olarak ayrılması gözlenmektedir.



Şekil 1. 13. *Eğitim ve Test Verisi*

**Kaynak:** (Kumar, 2020)



## Sınıflandırma Performanslarının Değerlendirilme Yöntemleri

Birçok sınıflandırma algoritması bulunmasından dolayı hangisinin kullanılması gerektiğine karar vermek için bazı ölçütler bulunmaktadır. Ölçütler karışıklık tablosu (confusion matrix) kullanılarak hesaplanmaktadır. Bunlar; doğruluk (accuracy), duyarlılık (sensitivity), kesinlik (precision) ve F-ölçütü (F-Score)'dür (Gümüştas, 2019; Yaygın, 2019).

### Karışıklık Matrisi (Confusion Matrix)

Sınıflandırma algoritmaları kullanılarak yapılan çalışmalar modellerin performanslarını gözlemlemek amacıyla kullanılmaktadır. Karşılaştırma için gerçek değer ile tahmin değeri Tablo 1.1'den yararlanılarak tespit edilmektedir (Yaygın, 2019; Gümüştas, 2019). Tablo 1.1'de karışıklık tablosunun gerçek ve tahmin değerlerine göre yorumlanmasına yer verilmektedir.

Tablo 1. 1. *Karışıklık Matrisi*

		Tahmin Edilen Veri	
		Pozitif	Negatif
Gerçek Veri	Pozitif	<b>Doğru Pozitif (DP):</b> Gerçek veride pozitif olan, tahmin edilen veride de pozitif olarak sınıflandırılmıştır.	<b>Yanlış Negatif (YN):</b> Gerçek veride pozitif olan, tahmin edilen veride ise negatif olarak sınıflandırılmıştır.
	Negatif	<b>Yanlış Pozitif (YP):</b> Gerçek veride negatif olan, tahmin edilen veride ise pozitif olarak sınıflandırılmıştır.	<b>Doğru Negatif (DN):</b> Gerçek veride negatif olan, tahmin edilen veride de negatif olarak sınıflandırılmıştır.

### Doğruluk (Accuracy)

Doğruluk, doğru sınıflandırma sonuçlarının toplam gözlem sayısına bölünmesi ile elde edilmektedir. Doğruluk ile ilgili eşitlik;

$$\text{Doğruluk} = \frac{DP + DN}{DP + DN + YP + YN}$$

şeklindedir. Eğer veri setinin dağılımı dengeli değilse doğruluk oranı veri kümesinin yoğun olduğu tarafa kayarak hatalı bir sonuç elde etmektedir (Yaygın, 2019; Gümüştas, 2019).

### **Kesinlik (Precision)**

Kesinlik, sınıflandırılması ve gerçek değeri pozitif olan gözlem değerlerinin, pozitif tahmin değerlerine oranlanması ile elde edilmektedir (Yaygın, 2019; Gümüştas, 2019). Kesinlik ile ilgili eşitlik;

$$Kesinlik = \frac{DP}{DP + YP}$$

şeklindedir.

### **Duyarlılık (Sensitivity)**

Duyarlılık, sınıflandırılması ve gerçek değeri pozitif olan gözlem değerlerinin, pozitif ve negatif tahmin değerlerine oranlanması ile elde edilmektedir. Duyarlılık ile ilgili eşitlik;

$$Duyarlılık = \frac{DP}{DP + YN}$$

şeklindedir. Sınıflandırılması ve gerçek değeri pozitif olan gözlem sayısı arttıkça duyarlılık artmaktadır (Yaygın, 2019; Gümüştas, 2019).

### **F-Ölçütü (F-Score)**

Kesinlik ve duyarlılık değerlerinin harmonik ortalaması ile elde edilmektedir. 0 ile 1 arasında değerler almaktadır. F-Ölçütü ile ilgili eşitlik;

$$F \text{ ölçütü} = 2 \frac{Kesinlik * Duyarlılık}{Kesinlik + Duyarlılık}$$

şeklindedir. Sonucun 1'e yaklaşması ölçütün iyi olduğunu göstermektedir (Yaygın, 2019; Gümüştas, 2019).

### **Kappa Değeri**

Kappa değeri, iki ya da daha fazla değerlendiriciye ait sonuçlar arasındaki uyumun güvenilirliğini ölçmektedir. Bu değer -1 ile +1 arasında değerler almaktadır. Sonuç +1'e yaklaştıkça sonuçların birbiri ile uyumunun arttığını, -1'e yaklaştıkça sonuçların uyumsuzluğunun arttığını ve 0'a yaklaştıkça değer arasındaki uyumun tamamen rastlantısal olduğu gözlenmektedir (Yaygın, 2019; Gümüştas, 2019). Kappa değeri ile ilgili eşitlik;

$$Kappa = \frac{Po - Pe}{1 - Pe}$$

şeklindedir.

Kapa değeri, Tablo 1.2’de yer alan bilgiler doğrultusunda değerlendirilmektedir (Kılıç, 2015).

Tablo 1. 2. *Kappa Değeri ve Yorumu*

<b>Kappa Değeri</b>	<b>Yorum</b>
< 0	Şansa bağlı gelişen kötü uyum
0.01 – 0.20	Önemsiz düzeyde uyum
0.21 – 0.40	Zayıf düzeyde uyum
0.41 – 0.60	Orta düzeyde uyum
0.61 – 0.80	İyi düzeyde uyum
0.81 – 1.00	Çok iyi düzeyde uyum

#### **1.1.4.2. Makine Öğrenmesinin Eğitimde Kullanımı**

Ölçme, bir değişkenin gözlemlenerek, gözlem sonuçlarının sayı ve sembollerle ifade edilmesi şeklinde tanımlanmaktadır (Turgut, 1995). Tan (2011) ölçmeyi, değişkenlerin sayı ya da sembol ile ifade edilmesi süreci olarak kabul etmektedir. Linn ve Gronlund (2000) bireyin sahip olduğu özelliğin sayısal olarak betimlenmesi şeklinde tanımlamaktadır. McMillan (2008) için ölçme, bir nesnenin miktarının belirlenmesi için kullanılan işlem olarak değerlendirmektedir. Eğitimde ölçme ve değerlendirme ise, öğrencilerde meydana gelen davranış değişikliklerinin ortaya çıkarılması, öğrenci ile ilgili alınacak kararların belirlenmesi, öğrencinin durumunun sayısal olarak ele alınması şeklinde ifade edilmektedir (Toptaş ve Şen, 2021).

Eğitim alanında veri ile ilgili işlemler, bilgi ve iletişim teknolojileri (BİT) aracılığıyla gerçekleştirilmektedir. Son zamanlarda veri hacmindeki artış ile eğitim alanı karar verme sürecinde farklı analitik yapıya gereksinim duymaktadır. Bu nedenle verilerin depolanması, analiz edilmesi ve anlamlandırılması süreçlerini kapsayan büyük veri kavramı eğitim alanında de rol almaya başlamıştır (Toptaş ve Şen, 2021). Günümüzde eğitim kademelerinin her basamağında veriler depolanmaktadır. Bu durum eğitim alanında gelenekselden endeksli yaklaşıma geçişin ilk adımları olarak nitelendirilmektedir. Genel kullanımda büyük veri

bireysel verilerden kitlesel verilere ulaşmaktadır. Eğitim alanında ise, kitlesel verilerden bireysel verilere geçiş söz konusudur (Toptaş, 2021).

Araştırmacılar, büyük verinin eğitim alanında önemli bir yer alacağı öngörmektedir. Eğitimin kalitesinin artırılması, bireyselleştirilmiş programların belirlenmesi, öğrenci akademik performansının artırılması, öğretmen yeterliliğinin ortaya çıkarılması, müfredatın planlanması, idari süreçlerdeki işlem basamaklarının iyileştirilmesi, ilgili kişiler tarafından öğrenci performansının kontrol altında tutulması gibi birçok şekilde veriler kullanılmaktadır (Toptaş ve Şen, 2021).

Ülkelerin eğitimlerinin incelenmesi amacıyla belirli dönemlerde uluslararası düzeyde sınavlar gerçekleştirilmektedir. Ekonomik İş Birliği ve Kalkınma Örgütü (OECD) tarafından üç yılda bir yapılan PISA bu sınavlardan birisidir. PISA ile öğrencilerin gündelik yaşam bilgi ve becerileri ölçülerek ülkelerin eğitim sistemlerinin durumu ortaya çıkarılmakta ve ülkeler birbirlerine göre karşılaştırılmaktadır (MEB, 2018). Bu tür sınavlar ile ülkeler kendi geçmiş ve gelecek sınavları hakkında yorum yapabilmektedirler. Bu noktada büyük veri, günümüz teknolojisine uygun hareket edebilecek güçte çalışabilmektedir (Toptaş ve Şen, 2021).

Eğitim alanında büyük veri mi kullanılmalı yoksa geleneksel yöntemler mi sürdürülmeli sorusu hala cevabı bulanamamış sorular arasında yer almaktadır. Büyük veri yapılandırıcı ve geleneksel yaklaşımın birleşimi ile yeni bir yaklaşım olarak kabul edilmektedir (Toptaş ve Şen, 2021). Geleneksel yaklaşımda; sonuca odaklanma, öğretmeni merkeze alma, tek yönlü bilgi aktarma, bilgilerin ezberlene bilirliliği, öğretimde “ne” sorusuna cevap aranması, parçadan bütüne anlayışı, bilginin sadece bilişsel boyutta olması, kısa zamanda çok iş yapılması, öğrenilenin sınavda sergilenme beklentisi ve öğrencinin not, diploma ya da sertifika ile ödüllendirilmesi yer almaktadır. Yapılandırıcı yaklaşımda ise; sürece odaklanma, öğrencinin merkeze alınması ve öğretmenin rehber görevini üstlenmesi, bilgilerin analiz, değerlendirme gibi süreçlerden geçirilmesi, öğretimde “neden, nasıl ve niçin” sorularına cevap aranması, öğrenci bilgiyi hangi yol ile kolay anlamlandıracaksa ona göre tümdengelim ya da tümevarım yöntemlerinin kullanılması, bilginin bilişsel, duyuşsal ve psiko-motor açıdan ele alınması, etkili öğrenmenin hedeflenmesi, öğrenilenin gerçek hayata uyarlanabilmesi ve öğrencinin proje ödevleri ile hayata hazırlanması yer almaktadır (Başol, 2015). Büyük veri eğitimde başarıyı “ne” getirir sorusu ile geleneksel yaklaşımı, analiz sürecindeki “nasıl” sorusu ile yapılandırıcı yaklaşımı kullanmaktadır (Toptaş ve Şen, 2021).

Geleneksel yaklaşımda yazılı sınavlar, çoktan seçmeli testler ile “neden, niçin” sorularına cevap aranmaktadır. Eğer bir sınıfta başarı yüksekse, başarının neden yüksek olduğu araştırılır. Büyük veri yaklaşımında ise, veriler bireysel olarak toplanarak kitlesel bir sonuca ulaşılmaya çalışılmaktadır. Burada bireysel alınan veriler analiz edildikten sonra gelecekteki performans ve davranış hakkında tahminde bulunarak buna göre planlama yapılmaktadır. Bu durum, eğitimde büyük verinin kullanılması ile bireysel programla ve materyal verimliliğine başarı getireceğini öngörmektedir. Diğer yandan büyük veri ile alternatif ölçme değerlendirme araçları kullanılarak öğrencilerin gerçek hayat ile bilgileri arasında ilişki oluşturacakları ödevler verilmektedir (Toptaş ve Şen, 2021).

Eğitimde verimlilik ile ilgili artış, verimliliğin sürdürülebilir ve nitelikli olması ile sağlanmaktadır. Sürdürülebilir öğrenme ise, öğrenci hakkında toplanan bireysel verilerin analiz edilmesi ve buna göre bir planlama yapılması ile sağlanmaktadır. Bireysel verilerden bir planlama yapılması büyük veri endeksli ölçme ve değerlendirme araçları ile gerçekleştirilmektedir (Brynjolfsson, Hitt ve Kim, 2011).

Siemens ve Long (2011) tarafından yayımlanan bir makalede eğitim alanı için birçok girişim bulunmakta ya da planlanmaktadır. Burada bilgi işlem cihazlarının kullanımının artırılması, esnek sınıf tasarımları ve yenilikçi görsel yapılar gibi yeni teknolojik gelişmelerden bahsedilmektedir. Yükseköğretimi etkileyecek en büyük gelişme büyük veri ve analitik anlayışının eğitime dâhil edilmesidir. Verilere dayalı kararlar kuramsal çıktıyı ve üretkenliği geliştirmektedir. Bu nedenle kararları verilere dayandırmak güçlü sonuçlar elde edilmesini sağlayacaktır (Brynjolfsson ve diğerleri, 2011). Sağlıkta, iş dünyasında sıklıkla kullanılan büyük veri müşteri ya da hasta hakkında beklenmedik sonuçlar ortaya çıkararak alanın gücünü arttırmaktadır (Axelrod, Brierley ve Vogel, 2011; Scism ve Maremont, 2010; Uzun ve Siemens, 2011). Eğitimde alanında ise geleneksel yöntemler kullanılmaktadır. Geleneksel yöntemler planlama ve kaynak belirleme süreçlerini, etkili öğrenme uygulamalarını, öğrenci profilleri gibi büyük miktarda veriyi kullanmakta başarısız olmaktadır. Eğitim için değişime gidilmesi ve öğrenmeye dayalı reform kararları almak için çalışmaların başlatılması önerilmektedir (Siemens ve Long, 2011).

Öğrenme ortamı karşılaştırılırken, mevcut okul değerlendirmelerinin çeşitli sınırlamalardan mustarip olduğu açıkça görülmektedir. Pedagojilerin çoğu öğrencilere çok az sayıda geri bildirim sağlar, öğretmenlerin rutin ödevlere not vermek için saatler harcamasını gerektirir, öğrencilere anlamanın nasıl geliştirileceğini gösterme konusunda öncelik sağlamaz ve öğrenme sürecini iyileştirebilecek dijital kaynaklardan yararlanmada başarısız

olmaktadır. Bu şekilde veriye dayalı yaklaşımlar öğrenmeyi, incelemeyi, öğrencilere ve öğretmenlere sistematik geri bildirim sunmayı mümkün kılmaktadır. Çalışmada veri madenciliği, veri analitiği ve web gösterge tabloları aracılığıyla gelişmiş araştırma, değerlendirme ve hesap verebilirlik potansiyeli incelenmiştir. Büyük veri, öğrenci performansı ve öğrenme yaklaşımları ile ilgili iç görüler için öğrenme bilgilerinin madenciliğini mümkün kılmaktadır (Manyika ve diğerleri, 2011). Bu yöntemler ile sadece test performansına güvenmek yerine, eğitmenler öğrencilerin ne bildiğini ve her bir öğrenci için hangi tekniklerin en etkili olduğunu analiz edebilmektedir. Öğretmenler, veri analitiğine odaklanarak öğrenmeyi çok daha detaylı şekillerde inceleyebilmektedirler (Castro, Vellido, Nebot ve Mugica, 2007). Çevrimiçi araçlar öğrencilerin okumaya ne kadar zaman ayırdıkları, elektronik kaynakları nereden aldıkları ve temel kavramlara ne kadar çabuk hâkim oldukları gibi çok daha geniş bir yelpazedeki öğrenci eylemlerinin değerlendirebilmektedir (West, 2012).

Büyük veri çerçevesinde ele alınan özelleştirilmiş ve dinamik öğrenme programı ile öğrencilerin öğrenme geçmişlerinden yola çıkarak, toplanan veriler ile öğrenciler için kişiye özel plan ve program hazırlanabilmektedir. Buna ek olarak, kariyer tahmini uygulamaları ile öğrencilerin güçlü ve zayıf yönleri belirlenebilmektedir (İntellipaat, 2016).

Büyük verinin eğitim alanına sağladığı önemli katkılardan birisi de yakın gelecek hakkında bir yordama oluşturulabilmesidir. Bahsi geçen yordama yapay zekâ, makine öğrenmesi ve veri madenciliği algoritmaları ile sağlanmaktadır. Bu durumun öğrenci performansının artırılması, öğretmen kalitesinin iyileştirilmesi gibi alana olumlu dönüşler sağlayacağı öngörülmektedir (Toptaş ve Şen, 2021).

Eğitimde büyük verinin kullanıldığı alanlar kısıtlıdır. Bahsi geçen alanda büyük verinin kullanılması eğitimcilerin eğitimdeki eksiklikleri görmeleri açısından önem arz etmektedir. Öğrencilerin akademik gelişimlerinin takip edilmesi, sınav sonuçlarının analiz edilmesi, ödevlerin takip edilmesi, müfredatın düzenlenmesi, alan ve meslek seçimi ile ilgili yönlendirme, eğitimdeki aksaklıkların tespit edilmesi ve çözümlenmesi gibi birçok alanda büyük veri analitiği kullanılmaktadır (Davenport, Jeanne ve Morison, 2010).

Toptaş ve Şen (2021)'e göre eğitimde büyük verinin kullanımı üç başlık altında ele alınmaktadır. İlki veri toplama ve depolama işlemlerinin gerçekleştirilebilmesi için uygun donanım ve yazılımların hazırlanmasıdır. İkincisi ise, eğitim alanında geleceğe yönelik tahmin çalışmaları gerçekleştirmek için analiz sürecinde Yapay Zekâ ve Makine Öğrenmesi algoritmalarından yararlanılmasıdır. Üçüncü ve en önemli madde ise, kitlesel sonuçlar ile

öğrenciye yaklaşmak yerine, öğrenciye özel bireyselleştirilmiş öğrenme metotlarının hazırlanmasıdır. Bireysel eğitim programları öğrencilerin kendi hızları ve algı düzeylerinde öğrenmeyi sağlayacağından akademik başarıyı da arttıracaktır. Son zamanlarda öğrenci kariyerinin erken tespit edilmesi ve bu doğrultuda olasılıklı modeller tasarlanması amacıyla öğrenci bilgi sistemleri büyük miktarda veriler ile çalışmak için eğitim alanında makine öğreniminin kullanılması gerektiğini öne sürmektedir (Kondakçı, Emil ve Beycioğlu, 2019).

Son yirmi yılda makine öğrenimi alanında önemli gelişmeler olmuştur. Bu alan bilgisayarla görme, konuşma, tanıma, doğal dil işleme, robot kontrolü ve diğer uygulamalar için pratik yazılım geliştirmek amacıyla tercih edilen yöntem olarak ortaya çıkmıştır (Jordan ve Mitchell, 2015). Makine öğreniminin eğitimi olumlu yönde etkileyebileceği birkaç alan vardır. Dijital Teknoloji ve Yönetim Merkezi tarafından yürütülen çalışmada eğitim verilerinin miktarındaki artışın eğitimde makine öğrenimi kullanımını desteklediği rapor edilmiştir (Mduma, Kalegele ve Machuve, 2019). Çeşitli okullar sınıflarda teknoloji kullanımı yoluyla kişiselleştirilmiş öğrenme deneyimlerini oluşturmaya başlamıştır. Ayrıca, çevrimiçi kurslar milyonlarca öğrenciyi kendine çekmiş ve öğrencilerin öğrenme çıktılarını iyileştirmeye ve toplanan verilerden yararlanmaya yönelik makine öğrenimi yöntemlerini uygulama ve geliştirme fırsatı sunmuştur (Lee, Chung ve Suh, 2017).

Toplanan veri miktarının artması ile öğrenme ve içerik analitiği, bilgi izleme, öğrenme materyali geliştirme ve erken uyarı sistemleri ile ilgili alanlar dâhil olmak üzere eğitim kalitesini artırmak için makine öğrenimi teknikleri kullanılmıştır (Mduma ve diğerleri, 2019). Bu tekniklerin eğitim amaçlı kullanımı, eğitim ortamlarında veri keşfetme yöntemlerinin kullanılması ve geliştirilmesinden dolayı umut verici bir alandır (Nunn, Avella, Kanai ve Kebritchi, 2016).

Eğitimde makine öğreniminin ilk uygulamalarından biri sınavların çoktan seçmeli testlerden, kısa yanıtlı cevaplara geçmesidir (Drabkin, 2017). Öğrencilerin serbest formdaki cevaplarının değerlendirilmesi Doğal Dil İşleme (NLP) ve Makine Öğrenimine dayanmaktadır. Otomatik puanlamanın etkinliği üzerine yapılan çeşitli çalışmalar, bazı durumlarda insan sınıflandırıcılardan daha iyi sonuçlar vermektedir. Ayrıca otomatik puanlama bir insandan daha hızlı puanlama sağlamakta ve bu da biçimlendirici değerlendirmede kullanıma yardımcı olmaktadır (Mduma ve diğerleri, 2019).

Kotsiantis (2012) tarafından yürütülen bir araştırmada, gelişmekte olan eğitimsel makine öğrenimi alanını tanımlayan yeni bir vaka çalışması yer almaktadır. Bu çalışmada, bir

öğrencinin gelecekteki performansını tahmin etmek için kullanılan bir makine öğrenimi regresyon yöntemi için veri seti olarak öğrencilerin temel demografik karakteristik verileri ve not verme verileri incelenmiştir. Benzer şekilde, bir öğrencinin okulu bırakma riskini tahmin etmek için eğitimciler, okullar ve politikacılar tarafından kullanılabilir bir tahmin modeli geliştirmeyi amaçlayan çeşitli projeler yürütülmüştür (Mduma ve diğerleri, 2019). Bu örneklerden yola çıkarak, öğretmenlerin okulu bırakma riski en yüksek olan öğrencileri belirlemelerine ve neden zorlandıklarını gözlemlemenin yanı sıra müdahalelere ilişkin iç görüş sağlamalarına yardımcı olabilecek bulut tabanlı öğrenme sistemlerini kullanan Akıllı Sınıflar tasarlanmıştır (Toptaş ve Şen, 2021).

Eğitimde makine öğrenimi uygulaması hala ele alınması gereken çeşitli zorluklarla karşı karşıyadır. Özellikle gelişmekte olan ülkelerde açık erişim veri kümeleri eksikliği vardır. Bunun dışında birçok araştırmacı değerlendirme prosedürlerinin ve ölçütlerinin okul yöneticileriyle ilgili olması gerektiği gerçeğini göz ardı etmektedir (Mduma ve diğerleri, 2019). Lakkaraju (2015) göre, değerlendirme sürecinde yalnızca yaygın olarak kullanılan makine öğrenimi algoritmalarına odaklanmak yerine, eğitimcilerin ihtiyaçlarını karşılayacak şekilde çalışmalar yürütülmelidir. Buna ek olarak; aynı çalışma, birçok çalışmanın sadece tahmine odaklandığını ortaya koymaktadır. Daha sağlam ve kapsamlı bir erken uyarı sistemi ile gelecekte risk altına girecek öğrencilerin belirlenmesi, öğrencileri okulu terk etme olasılıklarına göre sıralama ve risk altında olan öğrencilerin okulu terk etmeden önce belirlenmesi çalışmalarına öncelik verilmesi önerilmektedir. Bu nedenle gelişmekte olan ülkelerin, öğrencilerin okulu bırakmalarını için daha sağlam ve kapsamlı bir erken uyarı sistemi oluşturmaya odaklanması gerekmektedir. Ayrıca, sadece öğrenci düzeyindeki veri setlerine odaklanmak yerine okul düzeyindeki veri setlerine odaklanmaya ihtiyaç vardır. Bunun nedeni, okul bölgelerinin öğrencilere yardımcı olmak için genellikle sınırlı kaynaklara sahip olması ve bu kaynakların mevcudiyetinin zamana göre değişmesidir. Risk altındaki okulları belirlemek, yetkililerin riskten önce kaynak tahsisi için plan yapmasına yardımcı olacaktır. Ayrıca eğitim bağlamında, veri dengesizliği öğrenciyi elde tutma alanında çok yaygın rastlanılan bir sınıflandırma sorunudur. Çünkü kayıtlı öğrenci sayısının okulu bırakan öğrenci sayısına kıyasla daha fazla olması beklenmektedir (Thammasiri, Delen, Meesad ve Kasap, 2014).

Makine öğrenimi, yetkililerin sonuçları değiştiren önemli iç görüşler elde etmesine yardımcı olmak için daha iyi veriler oluşturmaya adım atılması sağlayabilmektedir. Öğrenciler eğitime devam etmek yerine okulu bıraktıklarında hem öğrenciler hem de topluluklar beceri, yetenek ve yenilikçilikten mahrum kalmaktadırlar. Öğrencilerin okulu bırakma sorununu ele



almak için gelişmiş ülkelerde kayıt, öğrenci performansı, cinsiyet ve sosyo-ekonomik demografi, okul altyapısı ve öğretmen becerileri ile ilgili ayrıntıları içeren karmaşık veri setlerini işlemek amacıyla tahmin modellerini bulmak için çeşitli tahmin modelleri geliştirilmiştir. Geliştirilen tahmine dayalı modellerin değerlendirilmeleri farklılık gösterme eğiliminde olmasına rağmen, yürütülen çalışmaların odak noktası gençlerin öğrenmeye devam etmesi için yatırım yapmak ve okul terklerini önlemek için en riskli öğrencilere müdahale etmektir. Bu nedenle, yöneticiler ve eğitimciler desteklenmektedir (Mduma ve diğerleri, 2019).

Toptaş ve Şen (2021) yayımladıkları makalede eğitim alanında büyük veri ile ilgili çalışmaları yönlendirebilecek ve çalışma alanı sağlayacak üç önemli konu belirtmiştir. Belirttikleri konulardan birisi; “*eğitim sektöründe Büyük Veri çalışmalarına altyapı oluşturması amacıyla hangi Yapay Zekâ (AI, Artificial Intelligent), makine öğrenmesi (ML- Machine Learning) programları ve algoritmalar kullanılabilir?*” sorusudur. Bu doğrultuda araştırmanın problemini, Türkiye Öğrenci Değerlendirmesi gerçek veri setinde öğretim elemanlarının sınıflandırılması için oluşturulacak Karar Ağacı, Rastgele Orman ve Yapay Sinir Ağları modellerinin incelenmesi oluşturmaktadır.

## **1.2.Araştırmanın Amacı ve Problemleri**

Eğitim alanında büyük veri çalışmalarına temel oluşturmak için makine öğrenmesi algoritmalarından hangilerinin alanda kullanılabileceğinin tespiti bu çalışmanın genel amacıdır.

Türkiye Öğrenci Değerlendirmesi gerçek veri seti ile öğretim elemanı kalitesinin araştırılması için öğretim elemanlarının performanslarının belirlenmesi ile ilişkisi olduğu düşünülen derse özel 28 soru ve 5 özellikten oluşan faktörler; Karar Ağacı, Rastgele Orman ve Yapay Sinir Ağları algoritmaları çerçevesinde incelenmiş ve bu üç algoritmanın sınıflandırma performansları araştırılmıştır. Bu amaç doğrultusunda belirlenen alt problemlere cevap aranmıştır:

1. Türkiye Öğrenci Değerlendirmesi veri setine göre öğretim elemanı başarısının belirlenmesinde Karar Ağacı, Rastgele Orman ve Yapay Sinir Ağları algoritmalarının sınıflandırma performansları;
  - 1.1.Doğru sınıflama oranına göre farklılaşmakta mıdır?
  - 1.2.Özgüllük oranına göre farklılaşmakta mıdır?
  - 1.3.Duyarlılık oranına göre farklılaşmakta mıdır?
  - 1.4.Kesinlik oranına göre farklılaşmakta mıdır?
  - 1.5.F1-İstatistiğine göre farklılaşmakta mıdır?

2. Karar Ağacı algoritmasına göre Türkiye Öğrenci Değerlendirmesi veri setinde öğretim elemanı sınıflamasının en önemli yordayıcıları nelerdir?
3. Rastgele Orman algoritmasına göre Türkiye Öğrenci Değerlendirmesi veri setinde öğretim elemanı sınıflamasının en önemli yordayıcıları nelerdir?
4. Yapay Sinir Ağı algoritmasına göre Türkiye Öğrenci Değerlendirmesi veri setinde öğretim elemanı sınıflamasının en önemli yordayıcıları nelerdir?
5. Öğretim elemanı sınıflamasının Karar Ağacı, Rastgele Orman ve Yapay Sinir Ağları yöntemlerine göre belirlenen manidar yordayıcıları farklılaşmakta mıdır?

### **1.3.Araştırmanın Önemi**

Büyük veri ve ona bağlı yöntemler sistematik ve düzenli bir şekilde analiz edildiği ve anlamlandırıldığı takdirde eğitim alanı içerisinde yer alan kişilere ve kurumlara katkı sağlayabilecek düzeydedir. Büyük veri ile makine öğrenmesi, veri madenciliği gibi birçok uygulamanın eğitim alanında kullanımı sağlanmaktadır. Bu doğrultuda bu çalışma, eğitim alanında makine öğrenmesi algoritmalarından hangilerinin kullanılması gerektiğini belirlemek amacıyla yürütülmüştür. Bu doğrultuda, gerçek veri seti üzerinden öğretim elemanı kalitesinin belirlenmesi amacıyla öğretim elemanı sınıflandırılması yapılmış ve uygun algoritma araştırılmıştır. Çalışma, eğitim alanında makine öğrenmesi algoritmaları kullanılarak öğrenci, öğretmen ve okul özelinde değerlendirmelerin yapılabilmesi için uygun algoritmanın araştırılması yönünden önemlidir. Ayrıca bu çalışma, eğitimde ölçme ve değerlendirme alanında gerekli varsayımların sağlanamadığı durumlarda analiz aşamasını kolaylaştırmak açısından da önem arz etmektedir. Makine öğrenmesinin eğitimde kullanılması ile ilgili çalışmaların literatürde bulunan sayısının kısıtlı olması nedeniyle, çalışmanın literatüre ve alana sağlayacağı katkı da bu çalışmanın bir başka önemi olarak görülmektedir.

### **1.4.Araştırmanın Varsayımları**

Bu çalışmada kullanılan 5820 Gazi Üniversitesi öğrencisinden toplanan Türkiye Öğrenci Değerlendirmesi veri setinin gerçek verilerden oluştuğu ve öğrencilerin uygulanan ölçüğe samimi yanıtlar verdiği kabul edilmektedir.

### **1.5.Araştırmanın Sınırlılıkları**

Bu araştırma, Türkiye Öğrenci Değerlendirmesi çalışmasına katılan 5820 Gazi Üniversitesi öğrencisinden elde edilen veriler, makine öğrenmesi sınıflandırma algoritmalarından Yapay Sinir Ağları, Rastgele Orman, C5.0 Karar Ağacı ve R-3.6.2 programı ile sınırlandırılmıştır.

## 1.6.Tanımlar

**Makine Öğrenmesi:** Makine öğrenimi, insanların öğrenme şeklini taklit etmek için veri ve algoritmaların kullanımına odaklanan bir yapay zekâ ve bilgisayar bilimi dalıdır.

**Karar Ağaçları:** Karar ağacı, sınıflandırma ve tahmin için en güçlü ve popüler araçtır. Her dalın testin bir sonucunu temsil ettiği ve her yaprak düğümün bir sınıf etiketi tuttuğu ağaç yapısı benzeri bir akış şemasıdır.

**Rastgele Orman:** Rastgele orman, tek bir sonuca ulaşmak için birden fazla karar ağacının çıktısını birleştiren, yaygın olarak kullanılan bir makine öğrenme algoritmasıdır

**Yapay Sinir Ağları:** Yapay sinir ağı, insan beyninin bilgiyi analiz etme ve işleme şeklini sembolize etmek için tasarlanmış bir bilgi işlem sisteminin parçasıdır.

**Türkiye Öğrenci Değerlendirmesi Veri Seti:** Bu veri seti, Ankara'daki (Türkiye) Gazi Üniversitesi öğrencileri tarafından sağlanan toplam 5820 değerlendirme puanını içermektedir. Toplam 28 kursa özel soru ve ek 5 özellik vardır.

## BÖLÜM II

### KAVRAMSAL ÇERÇEVE VE İLGİLİ ARAŞTIRMALAR

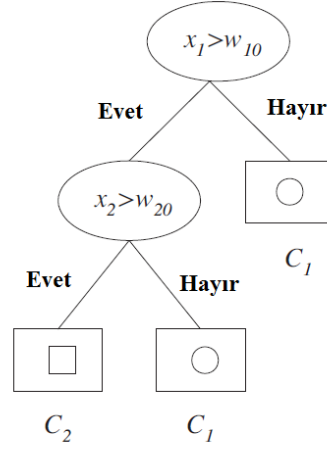
Bu bölümde ilk olarak makine öğrenmesi algoritmalarından Karar Ağaçları, Rastgele Orman ve Yapay Sinir Ağlarına değinilmiştir. Ardından bahsi geçen algoritmalar ile ilgili arařtırmalar eğitim alanında gösterdikleri performans ile sınırlandırılarak ele alınmıştır.

Algoritma, gerçek problemlere çözüm bulmak için tasarlanan yol şeklinde tanımlanmaktadır. Matematik ve Bilgisayar Bilimlerinin ortaklığından doğan sonlu işlemler kümesidir. Bilgisayar programlarında ağırlıklı olarak kullanılmakta ve bilgisayar dillerinin temel yapı taşı oluşturulmaktadır (Sitorus, 2015). Bu doğrultuda, matematik ve bilgisayarın ortaklığından doğan makine öğrenimi kapsamında regresyon ve sınıflandırma teknikleri için literatürde birçok algoritma yer almaktadır. Bunlardan sınıflandırma algoritmaları: Doğrusal Sınıflandırıcılardan Lojistik Regresyon (Logistic Regression) ve Naif Bayes Sınıflandırıcı (Naive Bayes Classifier), Destek Vektör Makineleri (Support Vector Machines), Karar Ağaçları (Decision Trees), Artan Ağaçlar (Boosted Trees), Rastgele Ormanlar (Random Forest), Nöral Ağlar (Neural Networks) ve En Yakın Komşu (Nearest Neighbor) şeklindedir (Alan ve Karabatak, 2020).

#### 2.1.Çalışmada Kullanılan Algoritmalar

##### 2.1.1. Karar Ağacı (Decision Tree)

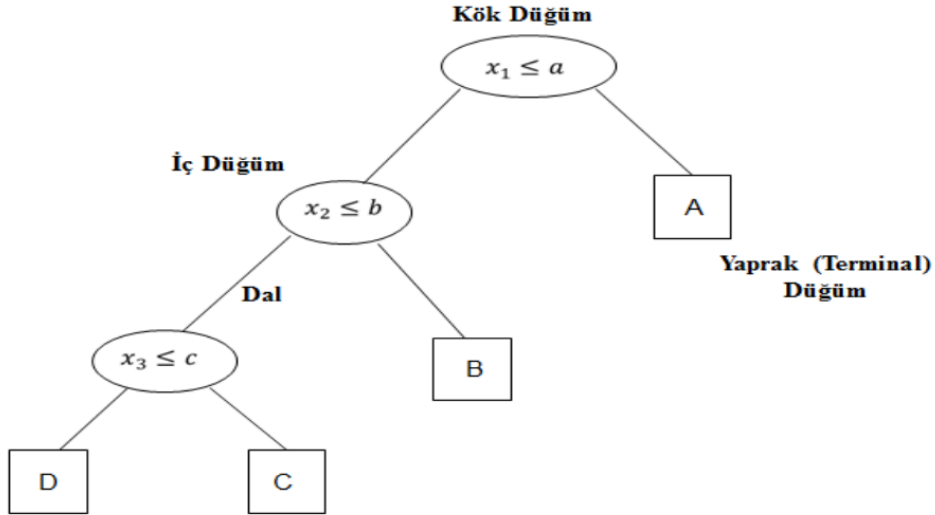
Karar ağaçları, eğitim kümesini kullanarak bir özellik doğrultusunda verileri art arda sınıflara ayıran denetimli öğrenme modelidir (Jankowski, Duch ve Grąbczewski, 2011). Karar ağacının amacı, girdi değişkenlerine dayalı bir hedefin değerini tahmin edecek model oluşturmaktır. Tahmini modelde, gözlem yoluyla belirlenen verilerin nitelikleri dallar tarafından temsil edilirken, verilerin hedef değerine ilişkin sonuçlarını yaprak temsil etmektedir (Boucheron ve Tagliaferri, 2019). Makine öğrenimi ve veri madenciliği çalışırken tahmine dayalı bir model tasarlanmaktadır. Modeller, verilerle ilgili gözlemleri ve verilerin hedef değerleriyle ilgili sonuçları eşleştirerek çalışmaktadır (Boucheron ve Tagliaferri, 2019).



Şekil 2. 1. Karar Ağacı

**Kaynak:** (Alpaydın, 2020)

Karar ağacı, örnek uzayının özyinelemeli bölümü olarak ifade edilen bir sınıflandırıcıdır. Algoritma kök, dal ve yaprakтан oluşan bir ağaç formundadır. Kök, ağacın ilk hücresidir, kökten çıkan dallar tahmin edilecek sınıfa giden yolu temsil etmektedir. Düğümler ise gözlemlerin sınıflandırılmasını sağlamaktadır. Düğüm, yaprak oluşana kadar devam edecek ve düğüm sayısı arttıkça ağaç karmaşık bir hale dönüşecektir. Yaprak son basamaktır ve sonucu bildirmektedir. Yaprak sınıflandırma basamağını temsil etmektedir (Maimon ve Rokach, 2007). Şekil 2.2’de algoritmanın kök, dal ve yapraklarının görseli almaktadır.



Şekil 2. 2. Karar Ağacı Kök, Dal ve Yaprak Dağılımı

**Kaynak:** (Yılmaz, 2014)

Ağacın yapısında bulunan dallar belirli bir kuralı temsil etmektedir. Her yeni dal oluşumunda sınıflandırma işlemi gerçekleşmemiş ise yeni bir karar düğümü oluşması beklenmektedir. Oluşan karar düğümlerinin sayısı derinlik olarak belirtilmektedir. Derinliğin sayısı ise veri setinin büyüklüğü ile homojen-heterojen olma durumuna bağlı kalmaktadır (Carvalho ve Freitas, 2004). Sistemde gerçekleşecek dallanmaların hangi ölçüt veya değişkene göre belirleneceği, karar ağaçları ile ilgili atılan en önemli adımlardan birisidir. Bu doğrultuda, testte kullanılmak üzere belirsizliği yüksek olan değişken kök düğümde kullanılmaktadır (Quinlan,1986).

Anlaşılması ve yorumlanması kolay, nicel ve kategorik verilerle çalışabilir, regresyon ve sınıflandırma teknikleri için uygun ve kayıp değere sahip veri setini işleyebilir olması karar ağaçlarının avantajlarıdır. Küçük değişikliklerden etkilenebilir yapıda olması, öğrenmedeki aşırılık, sürekli verilerle çalışırken başarısız olması, ağaç oluşturma ve budama sürecinde karmaşıklık yaşamaması ve sınıf sayısının fazla olduğu durumlarda başarısız model tasarlaması ise dezavantajlarıdır (Dixit, 2017).

Karar ağaçları, yüksek performans göstermesi ve verinin yapısı ile ilgili değerlendirme işlemlerini kolaylaştırmasından dolayı oldukça yaygın kullanılmaktadır. Algoritma kendi içinde birçok algoritmaya ayrılmaktadır. Bunlardan bazıları; AID, CHAID, CART, ID3, OEST, C4.5 ve C5.0'dir (Amasyalı, 2008).

## **Karar Ağaçları Algoritmaları**

### **AID Algoritması**

AID algoritması 1970'li yılların başında, Morgan ve Sonquist tarafından Michigan Üniversitesi (University of Michigan)'nde kullanılmaya başlayan, karar ağaçları için yazılmış ilk algoritmadır. Algoritma, bağımlı ve bağımsız değişkenler arasındaki ilişkinin incelenmesi amacıyla oluşturulmuş bir yazılımdır. Yeterince düzenli çalışmaması ve doğru ilişkiler kuramaması birçok araştırmacı tarafından belirtilmektedir (Akpınar, 2000).

### **CHAID Algoritması**

CHAID algoritması 1980 yılında G.V. Kass tarafından geliştirilen, AID algoritmasının geliştirilmiş halidir ve kategorik değişkenler arasındaki ilişkiyi en iyi açıklayacak şekilde bağımlı değişkenleri homojen alt gruplara bölerek ilerlemektedir. Kategorik değişkenleri kullandığı durum da ki-kareyi, sürekli değişkenleri kullandığı durum da ise F testini kullanmaktadır.

CHAID, CART algoritması ile benzerlik göstermektedir. Fakat CART algoritması ikili ağaçlar üretirken ilerlerken CHAID algoritması çoklu ağaçlanmalar ile ilerler. Bunu yaparken hem sürekli hem de kategorik değişkenleri kullanabilmesinden dolayı günümüzde popüler olarak kullanılan algoritmalarındandır (Pehlivan, 2006; Akpınar, 2000).

### **CART Algoritması**

CART algoritması 1984 yılında Breiman, Friedman, Olshen ve Stone tarafından geliştirilen, sınıflandırma ve regresyon algoritmasıdır. Veri setinde bağımlı değişkenlerin kategorik olduğu durumlarda sınıflandırma, sürekli olduğu durumlarda ise regresyon yapabildiği için sınıflandırma ve regresyon olarak adlandırılmaktadır. CART algoritması ikili ağaçlar halinde bölünerek ilerlemektedir. Bu süreçte kategorik değişkenler için Gini ve Twoing indeksini, sürekli değişkenler için en küçük kareler sapması kullanılmaktadır (Akpınar, 2000; Köktürk, 2012).

### **OUEST Algoritması**

OUEST Algoritması 1997 yılında Loh ve Shih tarafından geliştirilmiştir. Hızlı, etkin ve yansız olarak bilinen bir algoritmadır. Değişken seçimi ve düğüm noktalarının ayrımı CHAID ve CART'tan farklı olsada ikili ağaçlar halinde ilerlemesi açısından CART'a benzemektedir. Dallanma süreci ayrı ayrı ilerlemektedir (Pehlivan, 2006).

### **ID3 Algoritması**

1986 yılında Quinlan tarafından geliştirilmiştir. Karar ağacı algoritmaları arasında en basit algoritma olarak bilinmektedir. Sayısal değişkenler ya da kayıp değerler üzerinde herhangi bir budama yapılmaz. Bu yönüyle basit bir algoritma olarak kabul edilmiştir (Maimon ve Rokach, 2010).

### **C4.4 ve C5.0 Algoritması**

C4.5 algoritması, 1993 yılında Quinlan tarafından ID3 algoritmasının geliştirilmesi ile elde edilmiştir. Bölünme, örnek sayısının belirli bir sınıf altına düşmesi ile son bulur. Büyüme aşaması tamamlandıktan sonra hata budaması gerçekleşecektir. ID3'ten farklı olarak sayısal özellikleri işleyebilmektedir (Maimon ve Rokach, 2010).

### **C5.0 Algoritması**

C5.0 algoritması, C4.5 algoritmasının geliştirilmiş versiyonudur. Büyük verilere hizmet etmesi amacıyla geliştirilmiştir. C5.0 algoritması ile C4.5 algoritması birbirine benzemekle

birlikte birtakım farklılıklar da içermektedir. En temel farklılıkları C5.0 algoritması, C4.5 algoritmasına göre daha hızlı çalışmaktadır. Diğer algoritmalarından en büyük farkı ise, C5.0 normalizasyon yapmaktadır. C5.0 algoritması, dallanmanın her basamağında özellik kontrolü yapmakta, normalizasyona uğramış özellikler ile bilgi kazancı hesaplamakta, en iyi bilgi kazancına sahip özelliği ağaca eklemekte ve işleme basamaklarına budama ile devam etmektedir (Maimon ve Rokach, 2010).

C5.0 algoritması ile sonuçlar %90 oranında iyileştirilmiştir. İncelemeler neticesinde C5.0'ın, C4.5 algoritmasından 5,7 kez ve 240 kez daha hızlı olduğu ve daha kesin kurallar oluşturduğu sonucuna ulaşılmaktadır (Terlemez, 2008). C5.0 için dallanma tek bir düğüm ile başlamakta ve sınıflandırıcının belirlenmesi için entropi tabanlı bir bilgi kazancı kullanılmaktadır (Çakır, 2008)

C5.0 algoritması *kök düğümden yaprak düğüme uzanan karar kuralları* şeklinde tanımlanabilmektedir. Ağaç formunun yerine karar kurallarının daha kolay yorumlanabilir olmasından dolayı tercih edilmektedir. Karar kuralları öngörü kesinliğini arttırmak için budanabilir özelliktedir (Yakut, 2012).



### 2.1.2. Rastgele Orman (Random Forest)

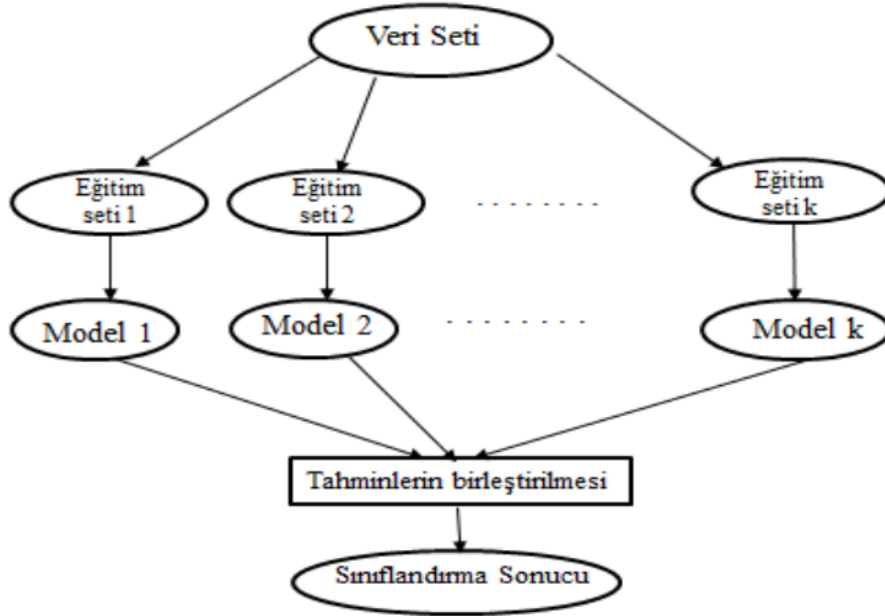
Rastgele Orman, 1998 yılında geliştirilen The Random Subspace tekniği ile Bagging yönteminin birleştirilerek, Breiman tarafından 2001 yılında geliştirilmiş topluluk öğrenme yöntemidir (Breiman, 2001). Rastgele Orman'ı diğer ağaçlardan ayıran özelliği içerisinde rastgelelik katmanını bulundurmasıdır. Burada her bir ağacın farklı bootstrap örneklemi bulunmakta ve bu, Sınıflandırma ve Regresyon Karar Ağaçlarının (CART) şeklini etkilemektedir. Karar ağaçlarında her düğümün tüm değişkenleri en iyi şekilde ayırması beklenirken; Rastgele Orman yönteminde her düğümün, düğümde kullanılacak rastgele seçilen tahmin edicilerden en iyisinin kullanılması beklenmektedir. Bu yöntemde ağaçlar bootstrap örnekleme ve düğüm ayrımlarında rastgele seçilen tahmin ediciler ile oluşturulmaktadır (Liaw ve Wiener, 2002).

Örneğin, bir çalışmada kullanılacak rastgele seçilmiş  $m$  tane girdi değişkeni ve bootstrap örnekleme ile seçilmiş her bir karar ağacı budanmadan en geniş hali ile bırakılmaktadır. Dikkat edilmesi gereken nokta,  $m$  tane tahmin edicinin toplam tahminci sayısından oldukça küçük olması gerektiğidir. Sınıflandırma işleminde ağaç sadece istenilen üyeler ile sınırlandırılır; regresyon da ise herhangi bir sınırlandırma olmadan yaprak düğümde az sayıda birim kalana kadar bölünme devam etmektedir (Yılmaz, 2014).

Rastgele Orman yöntemi, son zamanlarda popülerliğini arttırmış makine öğrenmesi yöntemidir. Diğer yöntemlerden farklı olarak tahmin doğruluğu ve model yorumlayabilme birlikteliği sağlamaktadır. Bu yöntemde kullanılan rastgele örnekleme ve topluluk stratejileri daha iyi genelleme yapabilmenin yanı sıra daha iyi ve doğru tahmin yapabilmeyi sağlamaktadır. Rastgele Orman'ın üç ana özelliği bulunmaktadır. Bunlar; farklı uygulamalara doğru tahminlerde bulunması, model eğitimi ile her özelliğin önemini ölçülebilmesi ve örnekler arasındaki yakınlığın model tarafından ölçülebilmesidir (Qi, 2012).

Rastgele Orman hem regresyon hem de sınıflandırma yöntemini doğal bir şekilde ele almaktadır. Eğitilmesi ve tahmin edilmesinin diğer yöntemlere göre daha hızlı olması, bir ya da iki parametreye bağlı hareket etmesi, genelleme hatası tahmini yapabilmesi, yüksek boyutlu problemlere doğrudan ve paralel olarak kolaylıkla uygulanabilmesi algoritmanın avantajları arasında yer almaktadır. İstatiksel olarak, rastgele orman yönteminin sağladığı ek özellikler ise; eksik değer ispatı, görselleştirme, aykırı değer tespiti ve denetimsiz öğrenmedir (Cutler, A., Cutler, D. R., ve Stevens, 2012). Ayrıca Rastgele Orman sezgisel strateji, ayrımcı analiz, vektör destek makineleri ve sinir ağları da dâhil olmak üzere birçok sınıflandırıcıyla karşılaştırıldığında daha iyi çalışmaktadır (Breiman, 2001).

Şekil 2.3'te yer alan görselde, rastgele orman algoritmasının birçok ağacının bir araya gelerek oluştuğu ve ardından bir sonuç değeri elde ettiği gözlenmektedir.

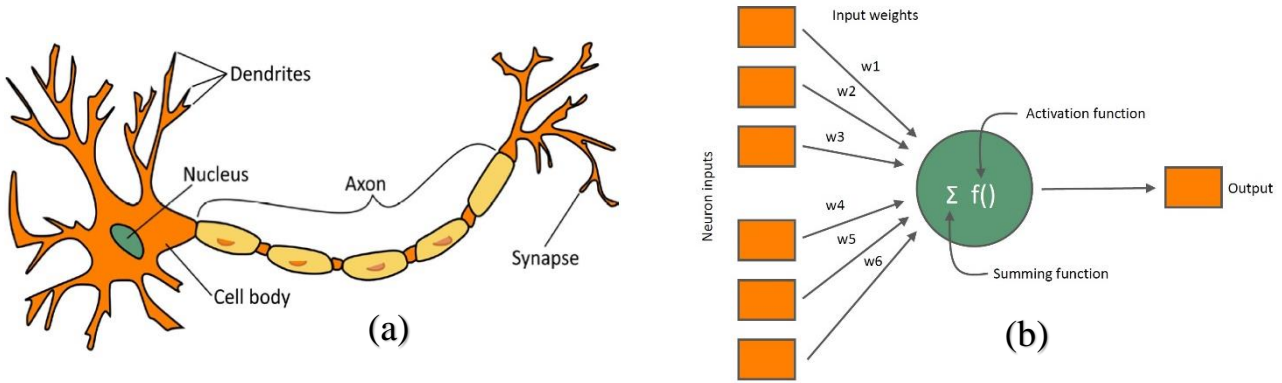


Şekil 2. 3. *Rastgele Orman*

**Kaynak:** (Yılmaz, 2014)

### 2.1.3. Yapay Sinir Ağları (Artificial Neural Networks)

Yapay sinir ağları, biyolojik organizmalarda öğrenme mekanizmasını sembolize eden popüler makine öğrenme teknikleridir. İnsan sinir sistemi, nöronlar olarak adlandırılan hücreleri içermektedir. Burada nöronlar, akson ve dentritlerin kullanımı ile birbirine bağlanmaktadır. Aksonlar ve dentritler arasındaki bağlantı bölgelerine sinaps adı verilmektedir. Bu bağlantılar Şekil 2.4 (a)'te yer almaktadır (Aggarwal, 2018).



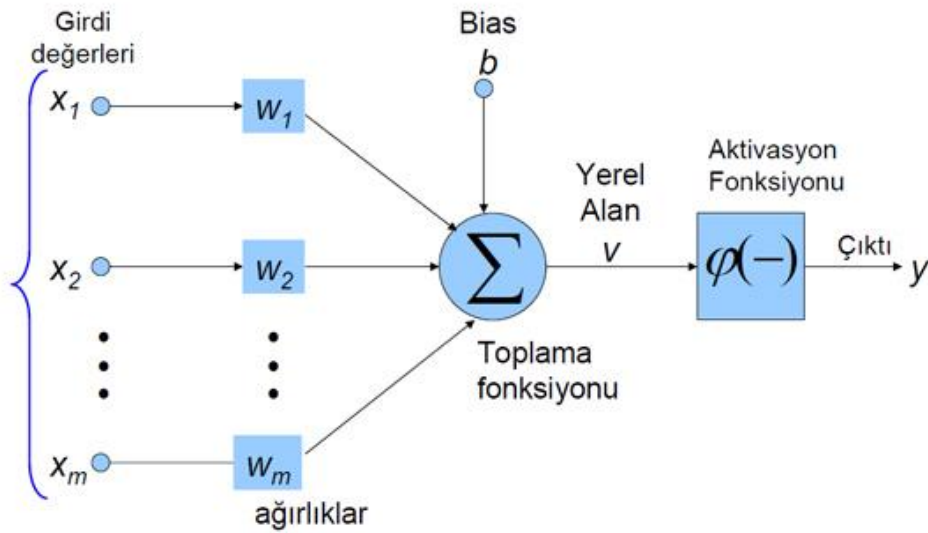
Şekil 2. 4. *Biyolojik Sinir Ağı (a) ve Yapay Sinir Ağı (b)*

**Kaynak:** (Ee Publishers, 2019)

Sinaptik bağlantıların gücü genellikle dış uyaranlara yanıt olarak değişmektedir. Bu değişim, öğrenmenin canlı organizmalarda nasıl gerçekleştiğidir. Bu biyolojik mekanizma, nöronlar olarak adlandırılan hesaplama birimlerini içeren yapay sinir ağlarında sembolize edilmektedir. Hesaplama birimleri, biyolojik organizmalardaki sinaptik bağlantıların gücüyle aynı rolü üstlenen ağırlıklar aracılığıyla birbirine bağlanmaktadır. Bir nörona her girdi, o birimde hesaplanan işlevi etkileyen bir ağırlıkla ölçeklenmektedir. Bu mimari Şekil 2.4 (b)'te yer almaktadır. Bir yapay sinir ağı, giriş nöronlarından hesaplanan değerleri çıkış nöronlarına yayarak ve ağırlıkları ara parametreler olarak kullanarak girdilerin bir fonksiyonunu hesaplamaktadır. Öğrenme, nöronları birbirine bağlayan ağırlıkların değiştirilmesiyle gerçekleşmektedir. Tıpkı biyolojik organizmalarda öğrenme için dış uyaranlara ihtiyaç duyulduğu gibi, yapay sinir ağlarında da dış uyaran, öğrenilecek fonksiyonun girdi-çıkış çiftlerinin örneklerini içeren eğitim verileriyle sağlanmaktadır (Aggarwal, 2018).

Yapay sinir ağları öğrenme ve karar sisteminin gelişmesini tetikleyen bilim dallarından birisidir. 1890 yılında beyin fonksiyonları ile ilgili yazılan ilk eserin yayınlanması ile başlayan YSA, 1970 öncesi ve sonrası diye ayrılmakta ve yakın gelecekte en önemli bilim dalı haline

geleceği ön görülen bir bilim dalıdır. YSA, insan beynine özgü; öğrenme, yeni bilgiler üretebilme ve oluşturabilme, keşfetme gibi yetenekleri herhangi bir destek olmaksızın otomatik olarak gerçekleştirmek için geliştirilmiş bir bilgisayar sistemidir. Geleneksel yöntemlerle bu yetenekleri uygulamak oldukça güçtür. Bu nedenle, yapay sinir ağları literatürde “uyarlayıcı (adaptif) bilgi işleme” olarak geçen beynin her yeni bilgiyi işlemesi ve işlevsel hale getirmesi durumunu karşılayan bir bilim dalı olarak kabul edilmektedir. Yapay sinir ağları, insan beynine özgü yetenekleri kullanarak çevreden gelen olaylara karşı tepki üreten bir sistem şeklinde tanımlanmaktadır. YSA; öğrenebilme, ilişkilendirebilme, sınıflandırabilme, genelleme, özellik belirleyebilme ve optimizasyon sağlayabilme yeteneklerini başarılı bir şekilde öğrenmekte ve uygulamaktadır. Tıpkı insan beyni gibi çalışan YSA, kendi “öz deneyim”ini kazanmakta ve bir sonraki süreçte aynı konu ile karşılaştığında kolaylıkla karar verebilmektedir. Birbirine hiyerarşik olarak bağlı ve paralel biçimde çalışan YSA’nın en temel görevi, kendisine gösterilen girdilere karşılık uygun çıktı sağlamaktır. YSA’nın bunu yapabilmesi için ağı ilgili örnekler verilerek öğrenme sağlanmakta ve sistem genelleme yapabilecek yeteneğe ulaştırılmaktadır. Böylece benzer olaylar ile karşılaşıldığında sistem çıktısı verisini kolaylıkla belirleyebilecektir (Öztemel, 2003).



Şekil 2. 5. Tek Katmanlı Tek Nöronlu Yapay Sinir Ağı

**Kaynak:** (Çınar, 2018)

Yapay sinir ağının temel bileşenini oluşturan yapı yapay sinir hücresidir. Şekil 2.5 incelendiğinde, YSA’nın bileşenleri; girdi değerleri ( $x$ ), sinapsların (diğer nöronlara bilgi

taşıyan özelleşmiş bağlantılar) ağırlığını gösteren ağırlık fonksiyonları ( $w$ ), toplama fonksiyonu ( $\Sigma$ ), aktivasyon fonksiyonu ( $f(-)$ ) ve  $y$  ile belirtilen çıktıdan oluştuğu gözlenmektedir (Rajapakse ve Omondi, 2006).

**Girdi (x):** Dışarıdan alınan hazır bilgiler şeklinde tanımlanmaktadır. Bilgi dışarıdan alınabildiği gibi sinir hücresinin kendisinden ya da başka hücreler tarafından bilgiye ulaşılabilir (Rajapakse ve Omondi, 2006).

**Ağırlık (w):** Girdi bilgisini önemini ve hücredeki etkisini gösteren yapı şeklinde tanımlanmaktadır. Burada değerler pozitif (+) ya da negatif (-) olabilir, bu ağırlığın etkisinin yönünü bildirmektedir. Fakat sıfır (0) olması ağ için önemlidir. Bu, ağırlığın hiçbir etkisinin olmadığını göstermektedir (Rajapakse ve Omondi, 2006).

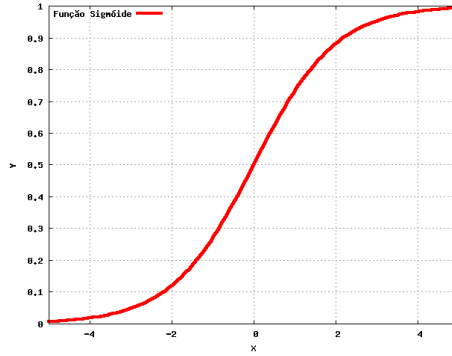
**Toplam Fonksiyon ( $\Sigma$ ):** Toplam fonksiyonu ile hücreye gelen net girdi değerleri hesaplanmaktadır. Bu aşamada farklı fonksiyonlar kullanılabilir. Net girdi değeri hesaplanırken, girdi değeri ile ağırlık çarpılmakta ve tüm değerler toplanmaktadır (Rajapakse ve Omondi, 2006).

$$NET = \sum_i^n x_i w_i$$

**Aktivasyon Fonksiyonu (f(-)):** İşlem basamağının son adımı olan aktivasyon fonksiyonu ile toplam fonksiyondan alınan girdi bilgisi işlenerek çıktı belirlenmektedir. Toplam fonksiyonunda olduğu gibi burada da farklı fonksiyonlar kullanılmaktadır. Her hücre için uygun fonksiyon diye bir kalıp söz konusu değildir, her problemin uygun fonksiyonu araştırmacı tarafından keşfedilmektedir (Öztemel, 2003). Aktivasyon fonksiyonu için en yaygın kullanılan formül sigmoid fonksiyonudur. Sigmoid fonksiyonu ile çıktı değerleri 0.0 ile 1.0 arasında sıkıştırılmaktadır. Bu özelliğinden dolayı yaygın olarak kullanılmaktadır. Şekil 2.6'da grafiği yer alan fonksiyonun matematiksel eşitliği;

$$f(x) = \frac{1}{1 + e^{-z}}$$

şeklindedir (Rajapakse ve Omondi, 2006). Matematiksel eşitliği verilen fonksiyonun grafiği şekil 2.6’da yer almaktadır.



Şekil 2. 6. *Sigmoid Fonksiyonu*

**Kaynak:** (Junior, 2011)

Aktivasyon fonksiyonu için farklı formüllerde kullanılmaktadır. Bunlardan bazılarını Şekil 2.7’de yer verilmiştir.

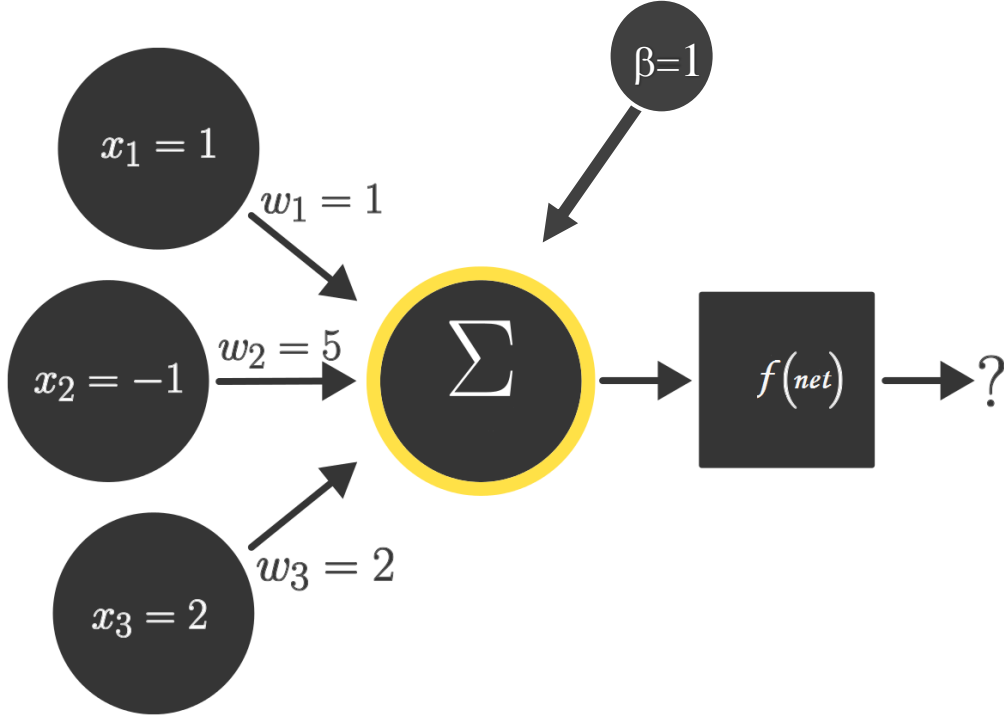
**Çıktı:** Aktivasyon fonksiyonu tarafından elde edilen bilgi şeklinde tanımlanmaktadır. Bu bilgi dış dünyaya gönderildiği gibi hücre kendisi için girdi olarak da kullanabilmektedir (Öztemel, 2003).

<b>Aktivasyon Fonksiyonu</b>	<b>Açıklama</b>
Lineer fonksiyon $F(s)=s$	Gelen girdiler olduğu gibi hücrenin çıktısı olarak kabul edilir.
Step fonksiyonu $F(s)= \begin{cases} 1 & \text{if } s > \text{eşik\_değer} \\ 0 & \text{if } s \leq \text{eşik\_değer} \end{cases}$	Gelen NET (s) girdi değerinin belirlenen bir eşik değerinin altında veya üstünde olmasına göre hücrenin çıktısı 1 veya 0 değerini alır.
Sinüs fonksiyonu $F(s)=\text{Sin}(s)$	Öğrenilmesi düşünülen olayların sinüs fonksiyonuna uygun dağılım gösterdiği durumlarda kullanılır.
Eşik değer fonksiyonu $F(s)= \begin{cases} 0 & \text{if } s \leq 0 \\ s & \text{if } 0 < s < 1 \\ 1 & \text{if } s \geq 1 \end{cases}$	Gelen bilgilerin 0 veya 1'den büyük veya küçük olmasına göre değerler alır. 0 ve 1 arasında değerler alabilir. Bunların dışında değerler alamaz.
Hiperbolik tanjant fonksiyonu $F(s)= (e^s + e^{-s}) / (e^s - e^{-s})$	Gelen NET girdi değerinin tanjant fonksiyonundan geçirilmesi ile hesaplanır

Şekil 2. 7. Aktivasyon Fonksiyonları

**Kaynak:** (Öztemel, 2003).

Bir yapay sinir hücresinin çalışma prensibi Şekil 2.8’de yer alan örnek soru üzerinden somutlaştırılmaya çalışılmaktadır.



Şekil 2. 8. Tek Katmanlı Tek Nöronlu YSA Örnek Soru

**Kaynak:** Yapay Sinir Ağları Örnek Sorular (2018)

$$y = \sigma(x_1w_1 + x_2w_2 + x_3w_3 + b)$$

$$\Sigma = (x_1w_1 + x_2w_2 + x_3w_3 + b) \quad (1)$$

$$\Sigma = (1.1 + (-1).5 + 2.2 + 1) \quad (2)$$

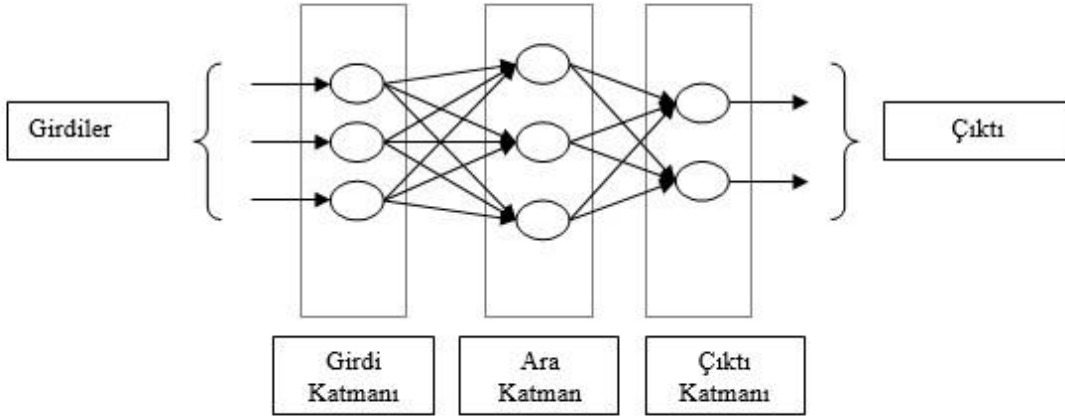
$$\Sigma = 1 \quad (3)$$

$$f(net) = \frac{1}{1+e^{-x}} \quad (4)$$

$$\text{Çıktı} = y = f(net) = 0,731 \quad (5)$$

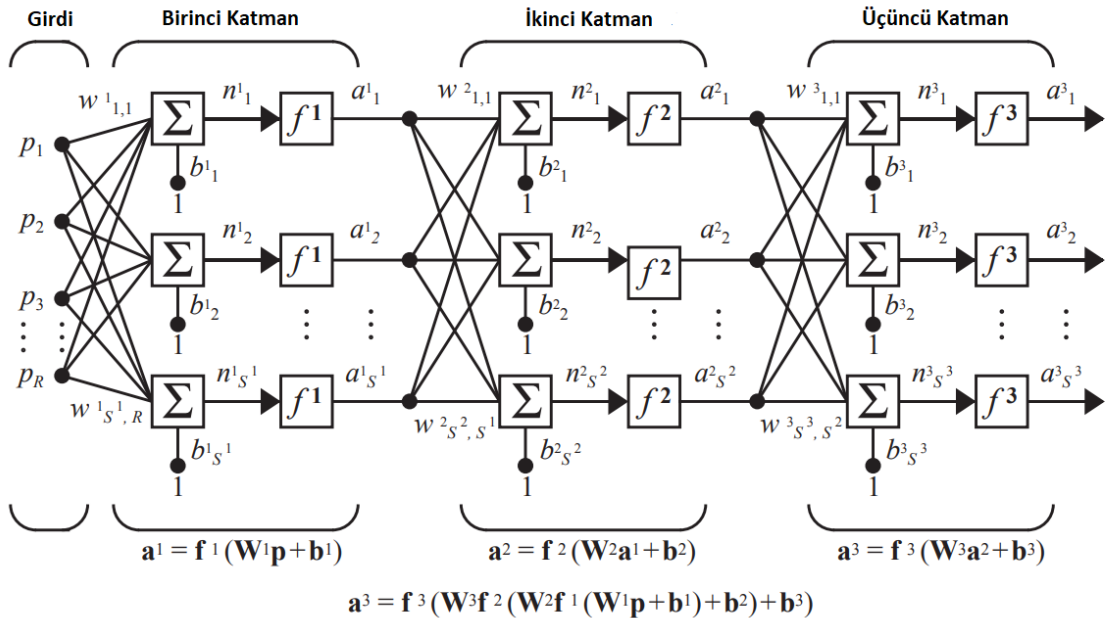


Yapay sinir ağının temel bileşeni yapay sinir hücresidir. Sinir hücreleri belirli bir kural ile bir araya gelerek sinir ağını oluşturmakta ve oluşan sinir ağı kendi içinde paralellik gösteren, Şekil 2.9'da yer alan görselde olduğu gibi üç katman (girdi katmanı, ara katman ve çıktı katmanı) halinde bulunan bir yapı şeklinde kullanılmaktadır (Öztemel, 2003). Şekil 2.9'da tek katmanlı üç nöronlu YSA modeli, Şekil 2.10'da modelin matematiksel adımları yer almaktadır.



Şekil 2. 9. Tek Katmanlı Çok Nöronlu Yapay Sinir Ağları

**Kaynak:** (Ağyar, 2016)



Şekil 2. 10. YSA Matematiksel Adımları

**Kaynak:** (Hagan, Demuth ve Beale, 1997)

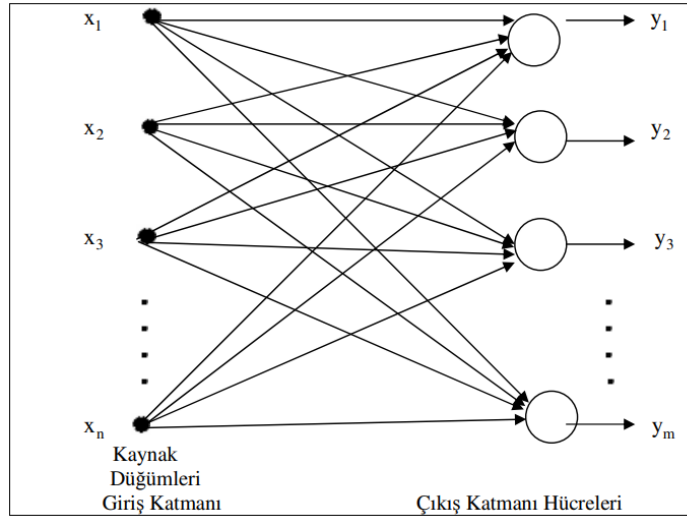
## Yapay Sinir Ağı Yapıları

### İleri Beslemeli Ağlar

Veri akışı, girdi katmanından çıktı katmanına doğru ilerlemektedir. Veriler kendi aralarında bir döngü ile ilerlemektedir. Yapay sinir ağlarının girdi katmanında yer alan veriler buradan ara katmana ve son olarak çıktı katmanına aktarılarak dış dünyaya geçiş yapmaktadır (Nasuhoğlu, 2019).

### Tek Katmanlı İleri Beslemeli Ağlar

Tek katmanlı ileri beslemeli sinir ağlarında birden fazla sinir hücresi bir araya gelerek bir katman oluşturulmakta ve girişten çıkışa doğru tek yönlü bir iletim sağlanmaktadır (Karakuzu, 1998). Şekil 2.11’de tek katmanlı ileri beslemeli ağın görseli yer almaktadır.

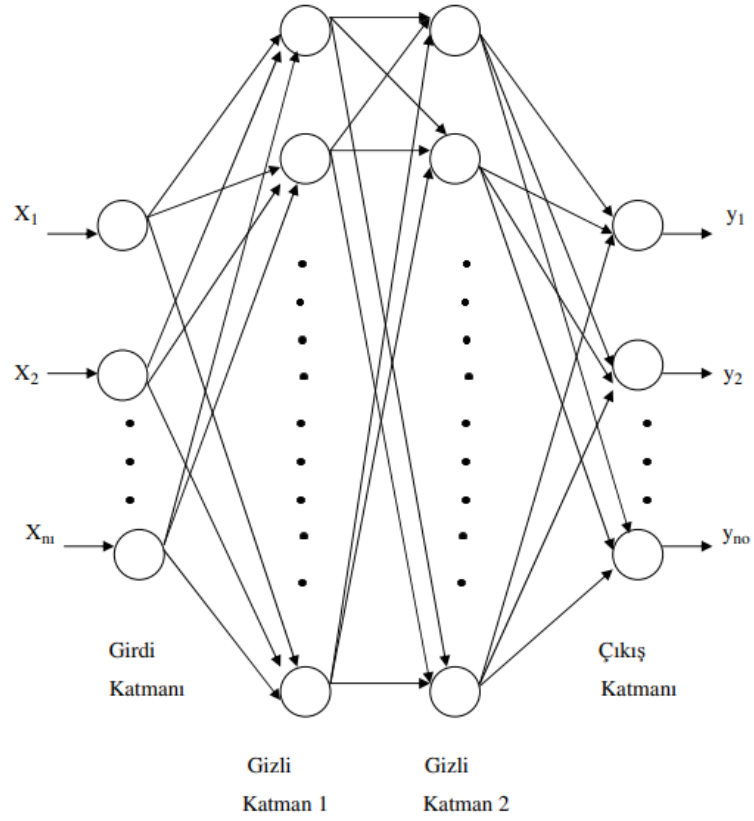


Şekil 2. 11. Tek Katmanlı İleri Beslemeli Ağlar

**Kaynak:** (Karakuzu,1998)

### Çok Katmanlı İleri Beslemeli Ağlar

Çalışma prensibi tek katmanlı ileri beslemeli ağlara benzemektedir. Çok katmanlı ağlarda, tek katmanlı ağlardan farklı olarak girdi ve çıktı katmanları arasında problem yapısına bağlı olarak oluşturulan çok sayıda gizli katmanın bulunmaktadır (Karakuzu, 1998). Şekil 2.12’de çok katmanlı ileri beslemeli ağın görseli yer almaktadır.

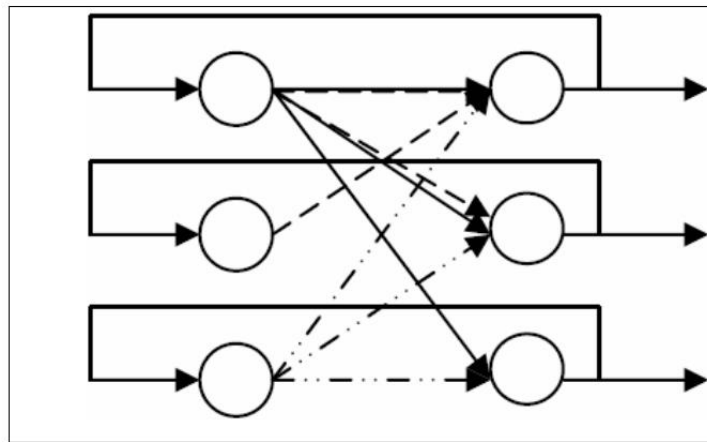


Şekil 2. 12. Çok Katmanlı İleri Beslemeli Ağlar

**Kaynak:** (Karakuzu,1998)

### Geri Beslemeli Yapay Sinir Ağları

Geri beslemeli yapay sinir ağları, hem ileri hem de geri yönde giriş bilgilerinin aktarılmasını sağlamaktadır. Girdi, ara katman ve çıktı olarak ilerleyen süreç; çıktı, ara katman ve girdi olarak da ilerleyebilmektedir (Asilkan ve Irmak,2009). Şekil 2.13'te geri beslemeli ağın görseli yer almaktadır.

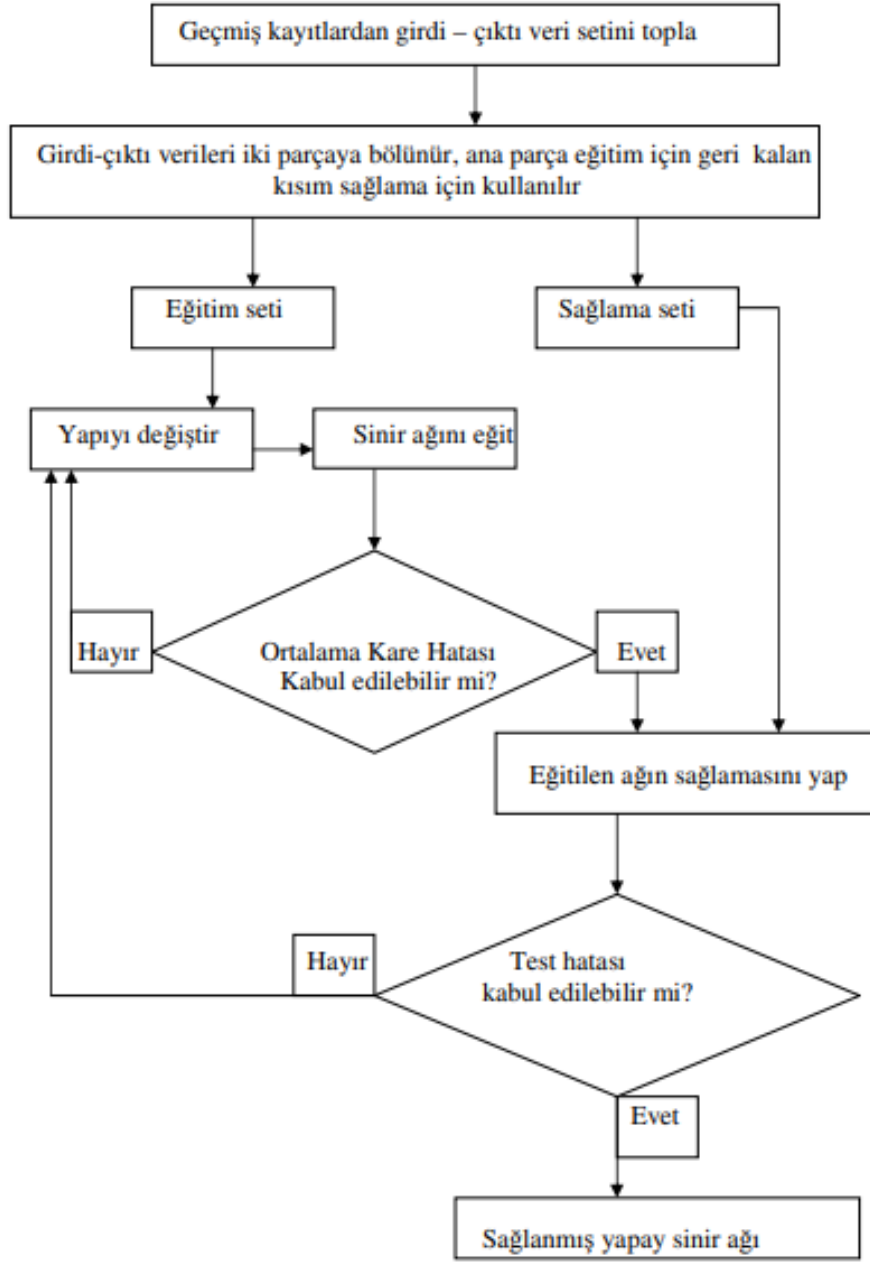


Şekil 2. 13. Geri Beslemeli Yapay Sinir Ağları

**Kaynak:** (Kaya, Oktay ve Engin, 2005)

## Yapay Sinir Ağlarında Eğitim,

Yapay sinir ağlarında eğitim sisteme girdi parametresi tanımlanmakla başlamaktadır. Burada parametreler işlevsel (operasyonel) parametreler (süre, yük, arıza modelleri, çevresel faktörler vs.) ve bakım parametreleri (onarım süreleri, yedek parça kullanılabilirliği vs.) içermektedir. Sisteme güvenilirlik, kullanılabilirlik ve sürdürülebilirlik gibi sistem performansını ölçen çıktı parametreleri tanımlanmaktadır. Girdi ve çıktılardan oluşan veri setinde, eğitim ve doğrulama için veri kümeleri hazırlanmaktadır. Genellikle verilerin üçte ikisi eğitim için kullanılırken, geri kalanı ise test ve doğrulama için kullanılmaktadır. Veri kümeleri, eğitim ve doğrulama için rastgele ayrılmaktadır. Yapay sinir ağları, eğitim setinin yardımıyla eğitilmektedir. Ortalama kare hatası minimuma ulaşana kadar eğitim gerçekleştirilmektedir. Eğitim basamağından sonra doğrulama veri kümesi ile YSA'nın geçerliği denetlenmektedir. Doğrulama aşamasında, yapay sinir ağının girdi setine doğru yanıt sağlama yeteneği değerlendirilmektedir. Burada yaygın olarak kullanılan strateji, eğitim sürecini tekrarlamak veya eğitim ve doğrulama sırasında ortalama kare hatasını gözlemlemektir. En iyi gizli katman sayısını önceden kestirmek zordur. Burada yapılması önerilen, tek bir gizli katman kullanmak ve her katmana az sayıda (örneğin, 2-5) sinir hücresi ile başlamaktır. Ardından eğitim ve doğrulama hataları gözlemlenmekte ve katman sayısı ile hücre sayısı arttırılabilmektedir. Son olarak, eğitilen ve doğrulanan yapay sinir ağı modelinin simülasyonu tasarlanmaktadır. Sistemin en iyi sonuca ulaşması için bu adımlar gerçekleştirilmektedir. Eğer çıktı istenilen şekilde değilse, en iyi sonuç elde edilene kadar adımlar tekrar etmektedir (Shishodia, Sekhon ve Rajpal, 2006). Şekil 2.14'te YSA algoritmasının eğitim aşamaları yer almaktadır.



Şekil 2. 14. Yapay Sinir Ağlarında Eğitim

**Kaynak:** (Karakuzu, 1998)

## Yapay Sinir Ağlarının Dezavantajları

Öztemel (2003)'e göre, yapay sinir ağlarının avantajlarının yanında birçok dezavantajları bulunmaktadır. Bilinen en önemli dezavantaj, YSA'nın donanım bağımlı çalışmasıdır. Paralel işlemciler ile çalışan ağlar, günümüzde çoğunlukla seri şekilde çalışan ve aynı zamanda tek bir bilgiyi işleyen makineler ile kullanıldıklarında bu durum zaman kaybına neden olmaktadır. Ayrıca ağın oluşturulması sürecinde belirli bir kuralın olmaması da bir başka dezavantaj olarak belirtilmektedir. Problem için en uygun ağ yapısı bulunurken deneme yanılma yolu kullanılmaktadır. Çünkü uygun ağın bulunamaması, problemin çözülememesine ya da düşük performanslı çözümlere yol açmaktadır. YSA kabul edilebilir düzeyde çözüm üreten bir sisteme sahiptir, bu nedenle en iyi çözümü sunamayacak ve garanti edemeyecektir. YSA'da katman sayısı gibi parametrelerin belirlenmesi için bir kural yoktur. Bu durum en iyi çözüm yolunu bulmayı güçleştirecek ve standart çözüm oluşturulmasını önleyecek önemli bir başka dezavantajdır. Yapay sinir ağının öğreneceği problemin ağa tanıtılması çok önemli bir adım olarak belirtilmektedir. Sadece numerik bilgiler ile çalışan YSA'da, problemin numerik gösterime dönüştürülmesi gerekmektedir. Uygun gösterim sisteminin kurulmamış olması problem çözümünü engellemekte ya da düşük performans göstermesine neden olmaktadır. YSA'da eğitimin ne zaman biteceğine karar vermek için geliştirilmiş bir yöntem yoktur. Süreçte hatanın belirli bir değerin altına indirilmesi eğitimin tamamlandığını göstermektedir. Fakat belirli bir kural dizisinin olmaması öğrenmenin en iyi düzeyde olduğunu göstermemektedir. En önemli dezavantaj ise, ağın davranışlarının açıklanamamasıdır. Bu durum ağın sonucuna olan güveni azaltmaktadır.

## 2.2.İlgili Araştırmalar

Bu kısımda makine öğrenmesi sınıflandırma algoritmalarından Karar Ağaçları, Rastgele Orman ve Yapay Sinir Ağları ile ilgili ulusal ve uluslararası literatürde yer alan, eğitim bilimleri alanındaki çalışmalar özetlenmiştir.

### 2.2.1. Uluslararası Araştırmalar

Drăgulescu, Bucos ve Vasiu (2015) tarafından çok sınıflı sınıflandırma probleminde ödev gönderimlerinin tahmin edilmesi üzerine bir çalışma yürütülmüştür. Öğrenci başarısızlığını tahmin etmek, eğitimcileri öğrenci performansını etkileyen faktörlere karşı koyma konusunda güçlendirebilecek önemli bir görev olarak görülmüştür. Bu doğrultuda çalışmada, öğrenci başarısızlığını tahmin etme gibi büyük bir problemin bir kısmına değinilmektedir. Bu sorunun çözümü için ilgili üniversitenin eğitim platformu tarafından toplanan gerçek veriler kullanılmıştır. Problem, üç olası sınıftan (çok sınıflı sınıflandırma) birini tahmin etmekten ibaret olduğundan uygun algoritmalar ve yöntemler seçilmiştir. Bu tahmin problemi ve kullanılan veri seti için en iyi yaklaşımı bulmak amacıyla birkaç deney yapılmıştır. Sonuç olarak, en iyi performans gösteren algoritmanın Rastgele Orman sınıflandırıcı olduğu tespit edilmiştir.

AL-Fakhry (2016) tarafından veri madenciliği tekniklerini kullanılarak C4.5 ve C5.0 algoritmalarının karşılaştırılması üzerine bir çalışma yürütülmüştür. Makalede veri madenciliği kavramının teorik konuları ve geliştirme adımları açıklanmaktadır. Kümeleme ve sınıflandırma süreci arasındaki farklar belirlenmektedir. C4.5 ve C5.0 sınıflandırıcılarının pratik yönleri ele alınmış ve karşılaştırmalı bir çalışma yapılmıştır. Her iki sınıflandırıcı için 30 hastaya test serumu uygulanmıştır. Sonuçlar olarak, yüksek çözünürlük sunmada C4.5 algoritmasının üstünlüğünü vurgulanmıştır.

Suh (2016) tarafından yayımlanan öğrenme analitiği ve eğitim veri madenciliği adlı çalışmada iki farklı veri seti için farklı algoritmalar ile model araştırılması yapılmıştır. Bunun için “Türkiye Öğrenci Değerlendirmesi Veri Seti” ve “Öğrenci Performansı Veri Seti” ayrı ayrı araştırılmıştır. Öğrencilerin dersi kaç kez aldıklarının araştırıldığı çalışmada, Türkiye Öğrenci Performansı veri seti için Naive Bayes, k-En Yakın Komşu, Lojistik Regresyon, J4.8 Karar Ağacı, JRip, Çok Katmanlı Algılayıcı ve ZeroR algoritmaları kullanılmıştır. En iyi performans gösteren algoritma %84.35 ile J4.8 Karar Ağacı olarak bulunmuştur.

Babić (2017) tarafından öğrencinin akademik motivasyonunu tahmin etmek için makine öğrenmesi yöntemleri üzerine bir çalışma yürütülmüştür. Akademik motivasyon, akademik performansla yakından ilişkilidir. Eğitimciler için, akademik motivasyonu yüksek öğrencileri tespit etmek kadar, akademik motivasyon eksikliği olan ergen öğrencileri tespit etmek de aynı derecede önemlidir. Araştırmacı, öğrencilerin Öğrenme Yönetim Sistemi (ÖYS)'nde yer alan derslerindeki davranışlarına dayalı olarak öğrencilerin akademik motivasyonunu tahmin etmek için bir sınıflandırma modeli geliştirmeye çalışırken, çalışma öngörülen öğrenci akademik motivasyonu ile ÖYS'deki davranışlar arasında bağlantılar kurmayı amaçlamaktadır. Bu araştırmaya Osijek'te yer alan Eğitim Fakültesi'ndeki öğrenciler katılmıştır. Üç makine öğrenimi sınıflandırıcısı (sinir ağları, karar ağaçları ve destek vektör makineleri) kullanılmıştır. Tüm sınıflandırıcılar başarılı performanslar göstermiştir, en iyi performans gösteren algoritma ise Sinir Ağları olarak bulunmuştur.

Luhaybi, Tucker ve Yousefi (2018) tarafından sınıflama yöntemleri kullanılarak öğrenci başarısızlığının tahmini üzerine bir vaka çalışması yürütülmüştür. Küreselleşen eğitim alanında, öğrenci performansını tahmin etmek, tahmine dayalı modelleri çok sayıda yönün etkilediği veri madenciliği ve makine öğrenimi araştırmacıları için merkezi bir konu haline gelmiştir. Çalışmada, öğrencilerin yükseköğretimdeki performansını değerlendirmek için sınıflandırma algoritmalarını uygulamaya ve üç ana nitelik kategorisinin bir kombinasyonuna dayalı olarak tahmin sürecini etkileyen temel özellikleri belirlemeye çalışılmıştır. Bunlar: kabul bilgileri, modülle ilgili veriler ve 1. sınıf final notlarıdır. Bu amaçla, 2015/16 akademik yılı için Londra Brunel Üniversitesi'nde bilgisayar bilimleri 2. sınıf öğrenci veri kümeleri üzerinde Karar Ağacı ve Naive Bayes sınıflandırma algoritmaları uygulanmıştır. Karar Ağacı algoritmasının en iyi performans gösterdiği sonucu ulaşılmıştır.

Hao, Galyardt, Barnes, Branch ve Wright (2018) tarafından bilgisayar eğitiminde etkisiz çevrimiçi öğrenci sorularının otomatik belirlenmesi adlı bir makale yayımlanmıştır. Bu araştırma ile etkisiz öğrenme sorularının otomatik olarak tanımlanması araştırılmaktadır. Etkisiz öğrenme sorularının anında ve doğru bir şekilde belirlenmesi, öğrencileri soruları gözden geçirmeleri için uyarmak ve uyarlanabilir soru revizyon önerileri sağlamak gibi büyük ölçekte olası otomatik kolaylaştırmaya kapı açmaktadır. Bunu başarmak için, Güneydoğu Amerika Birleşik Devletleri'ndeki büyük bir araştırma üniversitesinde üç dönem boyunca bir giriş programlama kursu tarafından uygulanan bir soru-cevap platformundan 983 soru toplanmıştır. Sorular önce manuel olarak üç hiyerarşik kategoriye ayrılmıştır: 1) öğrenmeyle ilgili olmayan sorular, 2) etkili öğrenmeyle ilgili sorular, 3) etkisiz



öğrenmeyle ilgili sorular. Manuel sınıflandırmanın değerlendiriciler arası güvenilirliği 0.88 olarak bulunmaktadır. Çok Terimli Naive Bayes, Lojistik Regresyon, Destek Vektör Makineleri ve Karar Ağacı dahil olmak üzere soruları otomatik olarak sınıflandırmak için dört farklı makine öğrenme algoritması kullanılmıştır. %90.1 ile NBM, %86 ile SVM, %84 ile DT ve %78.9 ile LR sonuçlarına ulaşılmaktadır

Krichevsky, Martynova ve Budagov (2019) tarafından yayımlanan yüksek lisans eğitim programında makine öğrenimi yöntemleri adlı makale, öğrencilerin çevrimiçi derslerde öğrenme deneyimlerini geliştirmek için erken tahminlerde bulunmaya yardımcı olacak bir süreç madenciliği yaklaşımı önermektedir. Bu tekniklerin etkinliğini ölçmek için süreç madenciliği özellikleriyle birlikte çeşitli makine öğrenme tekniklerinin etkisini araştırılmaktadır. Öğrenci verileri (değerlendirme notları, demografik bilgiler) ve olay günlüklerine dayalı haftalık etkileşim verileri (video ders etkileşimi, çözüm sunma süresi, haftalık harcanan zaman) makine öğrenmesi tekniklerine yönlendirmiştir. Çalışma, öğrencilerin performansının haftalık ilerlemesini izlemek ve genel performanslarını tahmin etmek için literatürde kullanılan dört makine öğrenimi sınıflandırma tekniğini (Lojistik Regresyon, Naive Bayes, Rastgele Orman ve K-En Yakın Komşu) değerlendirmektedir. Veri seti olarak, biri geleneksel özelliklere sahip ve diğeri süreç uygunluk testinden elde edilen özelliklere sahip iki veri seti kullanılmıştır. Çalışmada, dört tekniğin (LR, RF, Naive Bayes ve KNN) karşılaştırmalı bir analizi yapılmıştır. Sonuç olarak, kullanılan tekniklerin öğrencilerin performansını erken aşamada tahmin edebildiği gözlenmektedir. Süreç madenciliği özelliklerini geleneksel özelliklere entegre ederek bazı tekniklerin etkinliği artırılmıştır. LR ve Naive Bayes sınıflandırıcıları, sırasıyla veri seti bir ve veri seti iki için istatistiksel olarak anlamlı bir şekilde diğer tekniklerden daha iyi performans göstermiştir.

Otálora Orjuela (2019) tarafından yürütülen Kolombiya Üniversitesi'nin Ulusal Üniversitesi'nde de geçerli olan modeller adlı bir çalışmada, Kolombiya Üniversitesi'nin Uluslararası Ulusal Üniversitesi'ndeki lisans öğrencileri tarafından yaygın olarak seçilen derslerin ve müfredat yönergelerinin kaybıyla ilişkili kalıpları belirlemek için bir metodoloji geliştirilmiştir. Çalışma için, 2012-2017 yılları arasında bazı lisans müfredat programlarının gerçek akademik bilgileri üzerinde veri madenciliği ve makine öğrenimi yöntemleri kullanılmıştır. Çalışma planlarının müfredat esnekliğini örnekleyen, öğrenciler tarafından seçilen eğitim yollarındaki kalıplar belirlenmiştir. Bu anlamda, serbest seçim kredilerinin kaydı, öğrencilerin ait oldukları okullardan daha farklı okullarda yapıldığından sonuçlar dışsal

bir davranış göstermektedir. Sonuç olarak, tahmin ve sınıflandırma açısından algoritmaların performansları benzer sonuçlara sahip olduğu sonucuna ulaşılmaktadır.

Ortego ve Sánchez (2019) tarafından ilköğretim ve ortaöğretim öğrencilerinde okunan kitapların metin okunabilirlik derecesine göre sınıflandırılması için ilgili parametreler adlı çalışma yayımlanmıştır. Makalede, ilköğretim ve zorunlu ortaöğretimin ikinci ve üçüncü kademelerindeki öğrenciler için okuma kitabı seçimi söz konusu olduğunda okunabilirliğin en önemli parametrelerini belirlemeye çalışılmıştır. Veri madenciliği teknikleri aracılığıyla, okuma metinlerini, genel içerikleriyle ilgili olarak sözdizimsel, sözcüksel, anlambilimsel ve topolojik parametreler sayesinde otomatik olarak sınıflandırabilen bir hesaplama sistemi varsayılmıştır. Bu sürecin geçerliliği, farklı metinler arasında aynı doğrulama görevini yerine getirmek için kullanılan editör uzmanları tarafından dikkatli bir seçimle sağlanmıştır. Sonuç olarak Rastgele Orman en iyi performans gösteren algoritma olarak belirlenmiştir.

Yan ve Au (2019) tarafından makine öğrenimine dayalı çevrimiçi öğrenme davranış analizi üzerine bir çalışma yürütülmüştür. Makalenin amacı, öğrencilerin çevrimiçi öğrenme davranışları özellikleri ile ders notu arasında bir korelasyon analizi yapmak ve sınırlı verilere dayalı etkili bir tahmin modeli oluşturmaya çalışmaktır. Erişim günleri, kurs notuyla en yüksek korelasyona sahiptir, ardından isabet sayısı gelir ve bağlantı süresi öğrencilerin ders notuyla daha az ilgilidir. Öğrenci yaşı ve cinsiyeti, ders notu ile en düşük korelasyona sahiptir. İkili sınıflandırma modelleri, çok sınıflı sınıflandırma modellerinden çok daha yüksek tahmin doğruluğuna sahiptir. Sonuç olarak, Yapay Sinir Ağları en iyi performans gösteren algoritma olarak bulunmuştur.

Yekun ve Teklay (2019) tarafından optimum çok etiketli topluluk modeli ile öğrenci performans tahmini adlı bir çalışma yayımlanmıştır. Eğitim kalitesinin önemli ölçütlerinden biri öğrencilerin akademik performansıdır. Günümüzde eğitim kurumlarında öğrencilerin nasıl öğrendiğine ve veri madenciliği tekniklerini kullanarak performanslarını önceden nasıl geliştireceklerini keşfetmeye yardımcı olabilecek, öğrenciler hakkında bol miktarda veri depolanmaktadır. Bu çalışmada, lise öğrencilerinin gelecek dönem için beş derslik performansını tahmin eden bir öğrenci performans tahmin modeli geliştirilmiştir. Tahmin sistemi çok etiketli bir sınıflandırma görevi olarak modellenmiştir ve temel sınıflandırıcılar olarak Destek Vektör Makinesi (SVM), Rastgele Orman (RF), K-En Yakın Komşular (KNN) ve Çok Katmanlı Algılayıcıyı (MLP) kullanılmıştır. Sonuç olarak, Çok Katmanlı Algılayıcı en iyi performans gösteren algoritma olarak belirlenmiştir.

Altabrawee, Ali ve Ajmi (2019) tarafından makine öğrenimi tekniklerini kullanarak öğrencilerin performansını tahmin edilmesi üzerine bir çalışma yayımlanmıştır. Herhangi bir eğitim kurumunun nihai hedefi, öğrencilere en iyi eğitim deneyimini ve bilgisini sunmaktır. Bu hedefe ulaşmada ekstra desteğe ihtiyacı olan öğrencilerin belirlenmesi ve performanslarını artırmak için uygun aksiyonların alınması önemli bir rol oynamaktadır. Bu araştırmada, Al-Muthanna Üniversitesi (MU), College of Humanities tarafından sunulan bir bilgisayar bilimi dersinde öğrencilerin performansını tahmin edebilen bir sınıflandırıcı oluşturmak için dört makine öğrenme tekniği kullanılmıştır. Makine öğrenimi teknikleri arasında Yapay Sinir Ağı, Naive Bayes, Karar Ağacı ve Lojistik Regresyon bulunmaktadır. Bu araştırma, interneti bir öğrenme kaynağı olarak kullanmanın ve öğrencilerin sosyal ağlarda geçirdikleri zamanın öğrencilerin performansına etkisine ekstra önem vermektedir. Bu etkiler, öğrencinin öğrenmek için interneti kullanıp kullanmadığını ve öğrencilerin sosyal ağlarda geçirdikleri zamanı ölçen özellikler kullanılarak ortaya konmuştur. Modeller, ROC indeks performans ölçüsü ve sınıflandırma doğruluğu kullanılarak karşılaştırılmıştır. Ayrıca sınıflandırma hatası, kesinlik, geri çağırma ve F ölçüsü gibi farklı ölçüler hesaplanmıştır. Modelleri oluşturmak için kullanılan veri seti, öğrencilere verilen bir ankete ve öğrencilerin not defterine dayalı olarak toplanmıştır. Yapay Sinir Ağı (ileri beslemeli çok katmanlı ANN) modeli ile 0,807'ye eşit olan en iyi performans ve %77,04'e eşit olan en iyi sınıflandırma doğruluğu elde edilmiştir.

Achenie ve diğerleri (2020) tarafından küçük çocuklarda otizm taraması için makine öğrenimi stratejisinin kullanıldığı bir çalışma yürütülmüştür. Otizm Spektrum Bozukluğu taraması, erken tanı ve müdahale yoluyla prognozu iyileştirebilmektedir. Ancak zaman ve eğitim eksikliği pediatrik taramayı caydırabilmektedir. Yeni Yürümeye Başlayan Çocuklarda Otizm için Değiştirilmiş Kontrol Listesi, Revize (M-CHAT-R) yaygın olarak kullanılan bir eleme aracıdır, ancak takip eden sorular ve hataya açık insan puanlaması ve yorumlaması gerektirmektedir. Çalışmada OSB taramasının önündeki engellerin üstesinden gelmek için, özellikle İleri Beslemeli Yapay Sinir Ağını (fANN) kullanan otomatik bir makine öğrenimi yöntemini kullanılmıştır. Toplam örnek için, en iyi sonuçlar 18 madde kullanılarak %99.72 doğru sınıflandırma sağlamıştır. En iyi sonuçlar, Beyaz yürümeye başlayan çocuklar için 14 öge kullanılarak %99,92 doğru sınıflandırma ve Siyah yürümeye başlayan çocuklar için 18 öge kullanılarak %99,79 doğru sınıflandırma sağlanmıştır. Erkeklerde en iyi sonuçlar 18 madde kullanılarak %99.64 doğru sınıflandırma sağlarken, kızlarda 18 madde kullanılarak en iyi sonuçlar %99.95 oranında doğru sonuç vermiştir. Annenin eğitiminin 15 yıl veya daha az olduğu durumda (ön lisans ve altı) ve 16 madde kullanıldığında algoritma %99.75 doğru

sınıflandırmıştır. Anne eğitiminin 16 yıl veya daha fazla olduğunda (ön lisans derecesi üzerinde) ve 16 madde kullanıldığında algoritma %99.70 doğru sınıflandırma elde edilmiştir.

Musso, Rodríguez Hernández, ve Cascallar (2020) tarafından akademik yörüngelerde temel eğitim sonuçlarını tahmin etmek: bir makine öğrenimi yaklaşımı adlı makale yayımlanmıştır. Çalışmada geri yayılım algoritması ile çok katmanlı algılayıcı yapay sinir ağı modeli, özel bir üniversiteden 655 öğrenciden oluşan bir örnekleme not ortalaması, akademik kalıcılık ve derece tamamlama sonuçlarını sınıflandırmak için geliştirilmiştir. Bulgular, tüm sınıflandırmalar için yüksek düzeyde doğruluk göstermiştir. Yordayıcılar arasında öğrenme stratejilerinin not ortalamasının tahmin edilmesinde en büyük katkıya sahip olduğu gözlenmiştir. Başa çıkma stratejileri derecenin tamamlanması için en iyi tahmin edici ve üniversite programlarını bırakacak veya bırakmayacak öğrencilerin belirlenmesinde en büyük tahmin ağırlığına sahip olmaktadır. Sonuç olarak, %100 doğruluk oranı ile YSA mükemmel sonuç vermektedir.

DeVore, Yang, ve Stewart (2020) tarafından dengesiz fizik dersi sonuçlarını tahmin etmek için makine öğreniminin genişletilmesi üzerine bir çalışma yürütülmüştür. Makine öğrenimi algoritmaları, yakın zamanda öğrencileri bir fizik sınıfında A veya B alma olasılığı yüksek öğrenciler veya C, D veya F alma olasılığı yüksek öğrenciler olarak sınıflandırmak için kullanılmıştır. Bu çalışmada kullanılan performans ölçütleri, sonuç değişkeni büyük ölçüde dengesiz olduğunda güvenilir hale gelmiştir. Çalışma C, D ve F alacak öğrencilerin sınıflandırmasını genişletmeyi amaçlamaktadır. Bu çalışma için kullanılan örnekleme 7184 tane öğrenci oluşturmaktadır. Öğrencilerin yalnızca %12'si D veya F aldığında büyük ölçüde dengesizlik görülmüştür. Önceki çalışmayla aynı yöntemlerin uygulanması, D veya F durumlarının yalnızca %20'sini doğru sınıflandıran çok yanlış bir sınıflandırıcı üretilmiştir. Bu çalışmada rastgele orman makine öğrenmesi algoritması üzerinde durulacaktır. Rastgele orman karar eşiğini ayarlayarak, D veya F sonucunun doğru sınıflandırma oranını %46'ya yükseltmiştir.

Adnan, Habib, Ashraf, Shah ve Ali (2020) tarafından özellik ağırlıklarından yararlanılarak derin öğrenme teknikleriyle m-öğrencilerin performansını artırma adlı çalışma yayımlanmıştır. Mobil öğrenme (M-öğrenme), son on yılda eğitim ortamında büyük ilgi görmektedir. Etkili M-öğrenme için, mobil öğrencilerin (M-öğrenciler) tam gereksinimlerini tanımlayabilen verimli bir M-öğrenme modeli oluşturmak önemlidir. M-öğrenme modeli, M-öğrenenlerin mobil cihazlarla etkileşimi sırasında üretilen özelliklerden oluşmaktadır. Uyarlanabilir bir M-öğrenme modeli için, sadece öğrenme özellikleri gerekli

değildir, aynı zamanda bunların çeşitli M-öğrencileri, ağırlıkları ve karşılıklı ilişkileri için nasıl farklılık gösterdiğini belirlemek de önemlidir. Bu çalışma ile makine öğrenimi ve derin öğrenme tekniklerine dayanan sağlam ve uyarlanabilir bir M-öğrenme modeli önerilmektedir. Önerilen M-öğrenme modeli, M-öğrenenler için öğrenme özelliklerini, bunlara karşılık gelen ağırlıkları ve ilişkilendirmeyi dinamik olarak araştırmaktadır. Öğrenme özelliklerine dayalı olarak M-öğrenme modeli, M-öğrencilerini farklı performans gruplarına ayırmaktadır. M-öğrenme modeli daha sonra öğrenmeyi uyarlanabilir ve teşvik edici hale getirmek için M-öğrencilere uyarlanabilir içerik ve öneriler sunmaktadır. Karşılaştırmalı analiz için, beş temel makine öğrenimi modelinin tahmin doğruluğu, derin Yapay Sinir Ağı (derin YSA) ile karşılaştırılmıştır. Sonuçlar olarak, derin YSA ve Rastgele Orman modellerinin daha iyi tahmin doğruluğu sergilediğini sonucuna ulaşılmaktadır.

Richard ve Serrurier (2020) tarafından disleksi ve disgrafi tahminine makine öğrenimi yaklaşım adlı çalışma yürütülmüştür. Disgrafi, disleksi, dispraksi vb. gibi öğrenme güçlükleri öğrencilerin akademik başarılarını etkilemektedir. Ancak akademik zamanın ötesinde uzun vadeli sonuçları da bulunmaktadır. Dünya nüfusunun %5 ila %10'unun bu tür engellere maruz kaldığı yaygın olarak kabul edilmektedir. Erken çocukluk döneminde bu tür engelleri değerlendirmek için çocukların bir dizi test çözmesi gerekmektedir. Uzmanlar bu testleri puanlar ve çocukların notlarına göre belirli bir eğitim stratejisi gerektirip gerektirmediğine karar vermektedir. Değerlendirme uzun, maliyetli ve duygusal olarak acı verici olabilmektedir. Çalışmada, yapay zekânın bu değerlendirmeyi otomatikleştirmede nasıl yardımcı olabileceği araştırılmaktadır. Hem standart çocuklardan hem de disleksik ve/veya disgrafik çocuklardan el yazısı metin resimleri ve ses kayıtlarından oluşan bir veri seti toplamak, disleksik/disgrafik ve standart okuyucular/yazarlar arasındaki farkları analiz etmek ve bir model oluşturmak için makine öğrenmesi sınıflandırma teknikleri uygulanmaktadır. Algoritmalar model, resimler ve ses dosyaları analiz edilerek elde edilen basit özellikler üzerinde eğitilmiştir. Ön uygulamada kullanılan veri setine nispeten yüksek performanslar göstermektedir. Sonuç olarak, %74.7 ile Çoğunluk Sınıfı, % 90.8 ile Naive Bayes, %95.6 ile Lojistik Regresyon ve %96.2 ile Rastgele Orman elde edilmiştir.

Blasi ve Alsuwaiket (2020) tarafından Karar Ağacı ve YSA algoritmaları kullanılarak yükseköğretimde öğrenci suistimallerinin analizi adlı çalışma yayımlanmıştır. Yükseköğretim kurumlarının karşı karşıya olduğu önemli bir sorun öğrencilerin davranışlarının kötüye kullanılmasıdır. Bu çalışmanın amacı, üniversite kampüslerinde bunlara neden olan faktörleri belirleyerek bu suistimleri azaltmaktır. Sınıflandırma modelleri oluşturmak ve öğrencilerin

sınıflandırmak ve tahmin etmek için kurallar oluşturmak amacıyla J48 Karar Ağacı ve Yapay Sinir Ağları kullanılmıştır. Sonuçlar her iki sınıflandırma modeli için de yüksek sonuç elde edilmiştir.

Sobnath, Kaduk, Rehman ve Isiaq (2020) tarafından Birleşik Krallık'taki engelli öğrencilerin yükseköğrenim sonrası için makine öğrenimi yaklaşımı kullanılarak istihdam modelinin araştırıldığı bir çalışma yürütülmüştür. Çalışmada, makine öğrenimi ilkeleriyle büyük veri yaklaşımını kullanarak, mezuniyetten 6 ay sonra engelli öğrencilerin iş hayatına katılımıyla ilgili iyi bir tahmin edici oluşturmak için yararlı özelliklerden bir alt küme oluşturulmuştur. Diğerlerinin yanı sıra yaş, kurum, engellilik türü gibi özelliklerin önerilen istihdam modelinin temel belirleyicileri olduğu bilgisine ulaşılmıştır. Sonuç olarak, Karar Ağacı ve Lojistik Regresyon modellerinin, %96 doğrulukla en iyi sonuç sağladığı gözlenmiştir.

Zahour (2020) tarafından akademik ve mesleki rehberlik sorularının otomatik sınıflandırılması için makine öğrenimi yöntemleri ile karşılaştırmalı bir çalışma yürütülmüştür. Akademik ve mesleki rehberlik, giderek zorlaşan işgücü piyasasında önemli bir konudur. Bu bağlamda doğru karar verebilmek ve güçlü bir kariyer yolu oluşturmak için her öğrencinin ilgi alanı, ticareti, becerileri ve kişiliğini göz önünde bulundurmak çok önemlidir. Bu makale, dört makine öğrenimi algoritmasının sonuçlarının karşılaştırmalı bir çalışmasını sağlayarak eğitimsel ve mesleki rehberlik sorunsalını ele almaktadır. Kullanılan algoritmalar, okul oryantasyon sorularının otomatik olarak sınıflandırılmak ve John L. Holland'ın RIASEC Teorisi tipolojisine dayalı dört kategorize etmek için kullanılmaktadır. Sonuç Sinir Ağları'nın diğer üç algoritmadan daha iyi çalıştığını göstermektedir (Multiclass Decision Forest 0.75, Multiclass decision Jungle 0.75, Multiclass regression Logistic 0.79, Multiclass neural network 0.81).

Gorbanı ve Ghousi (2020) tarafından makine öğrenimi tekniklerini kullanarak öğrencilerin performansını tahmin etmede birbirinden farklı, yeniden örnekleme yöntemlerinin karşılaştırılması üzerine bir çalışma yürütülmüştür. Günümüz dünyasında teknolojinin ilerlemesi nedeniyle öğrencilerin performanslarını tahmin etmek en faydalı ve gerekli araştırma konuları arasında yer almaktadır. Veri madenciliği, özellikle öğrencilerin performansını analiz etmek için eğitim alanında son derece yararlıdır. Bu alandaki dengesiz veri kümeleri nedeniyle öğrencilerin performansını tahmin etmenin ciddi bir zorluk haline geldiği ve farklı yeniden örnekleme yöntemleri arasında herhangi bir karşılaştırmanın olmadığı bir gerçektir. Bu makale, iki farklı veri seti kullanarak öğrencilerin performansını tahmin ederken, dengesiz veri problemini ele almak için Borderline SMOTE, Random Over Sampler, SMOTE, SMOTE-ENN, SVM-SMOTE ve SMOTE-Tomek gibi çeşitli yeniden örnekleme tekniklerini

karşılaştırmaya çalışmaktadır. Ayrıca çok sınıflı ve ikili sınıflandırma arasındaki farklar ve özniteliklerin yapıları incelenmiştir. Dengesiz problemi çözmede yeniden örnekleme yöntemlerinin performansını daha iyi kontrol edebilmek için, çalışmada Random Forest, K-Nearest-Neighbor, Yapay Sinir Ağı, XG-boost, Destek Vektör Makinesi dâhil olmak üzere çeşitli makine öğrenimi sınıflandırıcılarını kullanılmıştır. Rastgele Orman sınıflandırıcısı, yeniden örnekleme yöntemi olarak SVM-SMOTE kullanıldığında diğer tüm modeller arasında en iyi sonucu elde etmiştir.

Kabathova ve Drlik (2021) tarafından farklı makine öğrenimi teknikleri kullanılarak üniversite öğrencilerinin dersi bırakmasını tahmin etmeye yönelik bir çalışma yürütülmüştür. Mevcut eğitim verilerine dayanarak öğrencilerin okulu terk etme durumlarını erken ve kesin olarak tahmin etmek, öğrenme analitiği araştırma alanının yaygın bir araştırma konusudur. Bu nedenle araştırmanın temel amacı, veri anlama, veri toplama aşamasının önemini vurgulamak, mevcut eğitim veri setlerinin sınırlamalarını vurgulamak, çeşitli makine öğrenimi sınıflandırıcılarının performansını karşılaştırmaktır. Çalışmada, e-öğrenme kursunda öğretmenler için mevcut olan sınırlı bir dizi özelliğin performans ölçütleri kapsamlı bir şekilde dikkate alınırca öğrencinin okulu bırakmasını yeterli doğrulukla tahmin edebileceği sonucuna ulaşılmaktadır. Bu doğrultuda, dört akademik yıldan toplanan veriler analiz edilmiştir. Lojistik Regresyon, Karar Ağacı, Rastgele Orman, Naive Bayes, Destek Vektör Makineleri ve Sinir Ağı kurs için seçilen algoritmalarıdır. Kesinlik, doğruluk ve F1 puanları karşılaştırıldığında en iyi performans tahmini %96.8 ile Rastgele Orman sınıflandırıcısı ile elde edilmiştir. %98.4 ile LR, %85 ile SVM ve %88.3 ile DT algoritmaları ile elde edilen sonuçlar RF ile karşılaştırıldığında küçük farklılıklarla ikinci en iyi sınıflandırıcı olarak belirtilmektedir. %76.5 ile Naive Bayes ve %88 ile Sinir Ağı modeli genel olarak en kötü sonuçlara sahip algoritmalar olarak kabul edilmektedir.

Latif, XianWen ve Wang (2021) tarafından üç seviyeli makine öğrenimi tekniğine dayalı öğrencilerin performansını tahmin etmek için akıllı karar destek sistemi yaklaşımı üzerine bir çalışma yürütülmektedir. Araştırmada Pakistanlı öğrencilerin akademik davranışlarını analiz etmek için kullanıcı dostu bir karar destek çerçevesi geliştirilmiştir. Bu makalenin amacı, üç seviyeli makine öğrenimi tekniğine dayalı bir akıllı karar destek sistemi (DSS) kullanan Pakistanlı öğrencilerin performansını analiz etmektir. Sinir ağı, Pakistanlı öğrenci başarısının tahmini için üç seviyeli bir sınıflandırıcı yaklaşımı kullanmaktadır. İngilizce ve Fizik derslerine katılan 1010 lisansüstü öğrencinin kendi kaydettiği bir veri seti kullanılmıştır. Bireysel öğrencinin algısını belirlemek için ön görüşme yapılmıştır. Anketin

istatistiksel anlamlılığını test etmek için ki-kare testi uygulanmıştır. İstatistiksel hesaplamalar ve verilerin hesaplanması, IBM SPSS sürüm 21.0'ın istatistiksel paketi kullanılarak yapılmıştır. Veri sınıflandırmasını iyileştirmek için yedi farklı algoritma test edilmiştir. Java tabanlı ortam, çok sayıda tahmin sınıflandırıcısının geliştirilmesi için kullanılmaktadır. Karar Ağacı algoritması en iyi doğruluğu gösterirken, Naive Bayes (NB) algoritması en az doğruluğu göstermektedir. Sonuçlar, sınıflandırıcının verimliliğinin %83,2'den %88,8'e kadar önerilen üç seviyeli bir şema kullanılarak iyileştirildiğini göstermektedir.

Mosquera Navarro (2021) tarafından yayımlanan Kolombiya'daki devlet okullarında ilk ve orta öğretim öğretmenlerinin psikososyal risk derecesini belirlemek için akıllı tekniklere dayalı sınıflandırma sistemi adlı doktora tezinde, Kolombiya'daki devlet okulu öğretmenleri arasında psikososyal risk tahminini iyileştirmek için akıllı bir algoritma geliştirmektedir. Model, sıvılarda yüzey geriliminin fiziksel teorisine bağlı yapay sinir ağı geri yayılımından oluşmaktadır. Bu çalışmanın amaçlarına ulaşmak için devlet okulu öğretmenleri, risk düzeyini belirlemek için meslek içi psikososyal risk faktörlerinin değerlendirilmesine yönelik bataryanın değerlendirilmesini gerçekleştirmiştir. Algoritmada psikososyal risk faktörlerini oluşturan değişkenler girdi, risk düzeyi ise çıktı olarak kullanılmıştır. Sonuç olarak, Geri Beslemeli Yapay Sinir Ağı %97,37 ile diğer algoritmalarından daha iyi bir sınıflandırma sağladığı sonucuna ulaşılmıştır.

### **2.2.2. Ulusal Araştırmalar**

Gündüz ve Fokoue (2015) tarafından yayımlanan öğrencilerin eğitimleri değerlendirmelerinde örüntü keşfi için istatistiksel veri madenciliği yaklaşımı adlı çalışmada Türkiye Öğrenci Değerlendirmesi veri seti ile eğitim sınıflandırması yapmak için sadece Rastgele Orman algoritması kullanılmış ve %98 performans gözlenmiştir.

Ahmed, Rizaner ve Ulusoy (2016) tarafından öğretim elemanı performansını tahmin etmek için veri madenciliğini yöntemlerinin araştırıldığı bir çalışma yürütülmüştür. Çalışma, öğretim elemanı performansını tahmin etmeye odaklanmakta ve eğitim sisteminin kalitesini iyileştirmek için öğrencilerin başarılarını etkileyen faktörleri araştırmaktadır. Türkiye Öğrenci Değerlendirme veri setinin kullanıldığı çalışmada J48 Karar Ağacı, Çok Katmanlı Algı, Naive Bayes ve Sıralı Minimal Optimizasyon gibi farklı veri sınıflandırıcıları üzerinde çalışmıştır. En iyi performans gösteren algoritma ise Karar Ağacı olarak belirlenmiştir.

Gök (2017) tarafından akademik başarının tahmin edilmesi sürecinde makine öğrenmesi algoritmalarının kullanılması üzerine bir çalışma yürütülmüştür. Çalışmada aile, demografik,



okul deęişkenleri kullanılarak 6,7 ve 8.sınıf öğrencilerinin akademik başarısını etkileyen faktörler ortaya konulmuştur. Öğrencilere 24 soruluk bir anket uygulanmıştır. Çalışmanın amacı, Türkçe ve Matematik dersleri ve dönem sonu başarı ortalamaları makine öğrenmesi denetimli öğrenme tekniklerinden sınıflandırma ve regresyon kullanılarak puan ve not tahmin edecek bir model tasarlamaktır. Sonuç olarak hem regresyon hem de sınıflandırma teknikleri ile öğrenci not tahmini için başarılı bir sonuç elde edilmiştir. Puan tahmin modelinde Rastgele Orman regresyonu ve not tahmini için Lojistik Sınıflandırıcı algoritmalarının en iyi performans gösterdiği sonucuna ulaşılmıştır.

Demirhan (2018) tarafından otizm spektrum bozukluęunun belirlenmesi için makine öğrenmesi algoritmaları kullanılarak bir çalışma yürütülmüştür. Çalışmada Destek Vektör Makineleri, K-En Yakın Komşu ve Rastgele Orman algoritmaları kullanılmıştır. Sonuç olarak, Rastgele Orman algoritması en iyi performans gösteren algoritma olarak belirlenmiştir.

Çifçi, Kaleli ve Günal (2018) tarafından yürütölen çalışmada makine öğrenmesi algoritmaları ile öğretim elemanı performansı tahmin edilmiş ve öznitelik araştırılmıştır. Çalışmada C4.5 Karar Ağacı, Naive Bayes, Derin Öğrenme ve K- En Yakın Komşu algoritmaları kullanılarak Türkiye Öğrenci Deęerlendirmesi veri seti ile öğretim elemanı performansının tahmin edilmesinde en iyi algoritma ve öznitelik araştırılmıştır. Çalışmada 2850 Gazi üniversitesi öğrencisinden toplanan veriler kullanılmıştır. Sonuç olarak, %97.70 ile Derin Öğrenme en iyi performans gösterirken hemen ikinci sırada %94,57 ile Karar Ağacı algoritması en iyi performans gösteren algoritma olarak bulunmuştur.

Kartal ve dięerleri (2019) tarafından öğrencilerin öğrenme stillerinin modellendięi bir araştırma yayımlanmıştır. Çalışmada C4.5 Karar Ağacı, Naive Bayes ve K-En Yakın Komşu algoritmaları araştırılmıştır. Sonuç olarak üç algoritmanın da yüksek performans gösterdiği sonucuna ulaşılmıştır.

Selvi (2020) tarafından makine öğrenmesi yöntemleri kullanılarak liseye geçiş sınavlarında öğrenci başarısının tahmin edilmesi üzerine bir çalışma yürütölmüştür. Çalışmada çok sınıflı makine öğrenmesi yöntemleri kullanmak için 4 farklı lisenin 9.sınıf ve okumaya devam eden 8.sınıf öğrencilerine 32 soruluk bir anket yapılarak gerekli veriler toplanmıştır. Elde edilen veriler ile en iyi modeli bulmak için bir araştırma yapılmıştır. Bunun için J48, PNN, Rastgele Orman, Karar Ağacı, RepTree ve Hoeffding Tree algoritmaları kullanılmıştır. Sonuç olarak, Rastgele Orman en iyi performans gösteren algoritma olarak bulunmuştur.

Osmanoğlu, Atak, Çağlar, Kayhan ve Can (2020) tarafından uzaktan eğitim ders materyalleri için duygu analizi üzerine bir makine öğrenmesi çalışması yürütülmüştür. Çalışmada, Karar Ağacı Sınıflandırıcı, MLP Sınıflandırıcı, XGB Sınıflandırıcı, Destek Vektör Sınıflandırıcı, Çok Terimli Lojistik Regresyon, Gaussian NB ve KNeighbors Sınıflandırıcı algoritmaları Python programlama dili ile modellenmiştir. Sonuç olarak, Lojistik regresyon algoritması ile 0.775 doğruluk oranı en iyi performans gösteren algoritma olarak belirlenmiştir.

Yıldız ve Börekçi (2020) tarafından makine öğrenmesi yöntemleri ile akademik başarı araştırılmıştır. Çalışmada Yapay Sinir Ağları, Karar Ağaçları, SVM, Rastgele Orman, K-En Yakın Komşu, Naive Bayes Lojistik Regresyon Sınıflandırıcı algoritmaları çerçevesinde araştırma yapılmıştır. Sonuç olarak, Yapay Sinir Ağları en iyi performans gösteren algoritma olarak bulunmuştur.

Afrin, Rahaman ve Hamilton (2020) tarafından makine öğrenmesi yöntemleri kullanılarak öğrenci memnuniyeti ortaya çıkartılmak istenilmiştir. Bunun için Türkiye Öğrenci Değerlendirme veri seti kullanılmış ve veri setinde öğrenci memnuniyetini belirleyen 10 etkili belirleyici ile kursla ilgili tahmin edicilerin 5 âdeti seçilmiştir. Çalışma için makine öğrenmesi yöntemlerinden Destek Vektör Makinesi, Çok Katmanlı Algılayıcı, Karar Ağacı, Karar Tablosu, k-En Yakın Komşu araştırılmıştır. Sonuç olarak makine öğrenimi yöntemlerinden, hem ders hem de öğretim elemanı ile ilgili faktörlerle eğitildiğinde, %80 ile %85 arasında doğrulukla öğrenci memnuniyetinin farklı perspektiflerini tahmin edebildiği ortaya konmuştur.

Şimşek ve Canbay (2021) tarafından covid-19 sürecinde uzak eğitim içi danışan gerekliliğinin makine öğrenmesi yaklaşımı ile belirlenmesi üzerine bir çalışma yürütülmüştür. Çalışmada K- En Yakın Komşu, Destek Vektör Makinesi, Naive Bayes, Karar Ağacı ve Rastgele Orman algoritmaları araştırılmıştır. Sonuç olarak, Destek Vektör Makinesi en iyi performans gösteren algoritma olarak belirlenmiştir.

Bu tez kapsamında ulusal ve uluslararası literatürdeki çalışmalara Karar Ağacı, Rastgele Orman ve Yapay Sinir Ağları ile eğitim veri seti sınırlılığında yer verilmiştir. Çalışmalar incelendiğinde, genel olarak Karar Ağacı ve Rastgele Orman algoritmalarının diğer algoritmalara görece daha iyi performans gösterdiği sonucuna ulaşılmıştır. Bu çalışma ile literatürde yer alan TÖD veri seti için uygun algoritma arayışına, farklı algoritma seçimi yapılarak katkı sağlanmıştır.

## BÖLÜM III

### YÖNTEM

Bu bölümde araştırma modeli, çalışma grubu verileri, verilerin analizi ve algoritmaların uygulaması üzerinde durulacaktır.

#### 3.1.Araştırmanın Modeli

Bu çalışmada Türkiye Öğrenci Değerlendirmesi gerçek veri seti kullanılarak öğretim elemanı sınıflandırılması için bir model tasarlanmıştır. Veri incelendiğinde hangi öğretim elemanının ne kadar olasılıkla seçildiği bilinmektedir. Fakat sınıflandırma çalışmaları bu işlemleri hızlandırmak ve gelecekteki veri setleri için bir alt yapı oluşturmak başka bir ifade ile hazır model oluşturmak amacıyla vardır (Shalev-Shwartz ve Ben-David, 2014). Bu doğrultuda, çalışmada mevcut durumun araştırılıp belirlenmesi hedeflediğinden araştırma modeli olarak betimsel model kullanılmıştır.

Lin (1976)'e göre, betimsel yöntem araştırılan durumun tam ve dikkatli bir şekilde tanımlanmaya çalışmaktadır. Araştırılan konu ya da durumun betimlenmesi, tasvir edilmesi ile sürdürülen araştırmalardır. Ele alınan örneklem ile ilgili bilgiler betimlenerek temel özellikler tasarlanmaktadır. Bu tür araştırmalarda tıpkı keşfedici araştırmalarda olduğu gibi “kim”, “ne” ve “neden” sorulara cevap aranmaktadır. Ancak betimleyici araştırmalar daha sistematik bir yol izlemektedir. Burada araştırılan konu ya da grup dışarıdan müdahale olmaksızın olduğu gibi doğal hali ile gözlemlenmekte ve betimlenmektedir. Betimsel yöntemler, incelenmek istenen araştırma konusu hakkında genel bir yorum yapabilmek, bakış açısı kazanabilmek için ideal bir yöntemdir. Nüfus sayımları betimleyici analizler için bir örnektir. Burada sayımın amacı, hem genel olarak ülke hakkında yorum yapabilmek hem de iller bazında değerlendirmeler yaparak kesin betimlemeler elde etmektir. Bu yöntemde neden sonuç ilişkisi ile ilgilenilmemektedir, ancak bu tür araştırmalarda da temel istatistik yöntemlerinden bazıları kullanılabilir. Bunlar; frekans dağılımı, ortalama değerler gibi istatistiksel yöntemlerdir. Betimleyici araştırmalar aynı konuda yapılacak açıklayıcı araştırmaların neye odaklanması gerektiği konusunda ipucu vermektedir. Örneğin, bir bölgede bir tür hastalığa yakalananları inceleyen bir betimsel araştırmanın ardından, hastalığın nedenleri hakkında açıklayıcı araştırma yürütülebilmektedir. Bu yöntemin avantajları, incelenmek istenen konunun kendi ortamında incelenmesidir. Dezavantajı ise, sonuçların ileri istatistiksel yöntemler ile analiz edilememesi ve araştırma sonuçlarının farklı yorumlanabilir olmasıdır.

### 3.2.Çalışma Grubunun Verileri

Bu çalışmada, Eğitimde Ölçme ve Değerlendirme alanı için önemli bir yazılım olan R programlama dili kullanılmıştır. Veri seti olarak Kaliforniya Üniversitesi Makine Öğrenmesi Veri Havuzunda (University of California, Irvine, USA) yer alan “Türkiye Öğrenci Değerlendirmesi Veri Seti (Turkey Student Evaluation Data Set)” adlı gerçek veri seti (Dosya: turkiye-student-evaluation\_R\_Specific.csv) kullanılmıştır. Veri seti, üç öğretim elemanı tarafından toplam 13 ders alan 5820 Gazi Üniversitesi öğrencisinin öğretim elemanı değerlendirme sonuçlarını içermektedir. Çalışma için 28 adet derse özel soru ve ek 5 özelliğten oluşan 5’li likert tipi ölçek geliştirilmiştir. Ölçekte “*Kesinlikle Katılmıyorum, Katılmıyorum, Kararsızım, Katılıyorum, Kesinlikle Katılıyorum*” ifadelerine yer verilmiştir. İfadelerin sayısal karşılığı (1,2,3,4,5) şeklinde belirlenmiştir (Gunduz ve Fokoue, 2013). Tablo 3.1’de ölçekte yer alan toplam 33 adet niteliğe ait özelliklere yer verilmiştir.

Tablo 3. 1. *Türkiye Öğrenci Değerlendirmesi Veri Seti*

Özelliğın Adı	Açıklama	Değerler
<b>egitmen</b>	Eğitmen tanımlayıcısı	{1,2,3}
<b>sınıf</b>	Alınan dersin ismi	{1-13}
<b>tekrar</b>	Öğrencinin dersi alma sayısı	{1,2,3}
<b>katılım</b>	Öğrencinin derse katılım seviyesi	{0,1,2,3}
<b>zorluk</b>	Öğrencinin yorumuna göre dersin zorluk derecesi	{1,2,3,4,5}
<b>Madde 1</b>	Dönem başında bildirilen ders içeriği ve öğretim metodu ve kıymetlendirme sistemi.	{1,2,3,4,5}
<b>Madde 2</b>	Dönem başında kursun amaç ve kazanımları açıkça bildirilmiştir.	{1,2,3,4,5}
<b>Madde 3</b>	Kurs, kendisine atanan kredi düzeyine karşılık verebilecek niteliktedir.	{1,2,3,4,5}
<b>Madde 4</b>	Kurs, ders programında belirtildiği şekilde işlendi.	{1,2,3,4,5}
<b>Madde 5</b>	Sınıf içi tartışmalar, verilen ödevler, yapılan uygulama ve çalışmalar yeterli düzeydeydi.	{1,2,3,4,5}
<b>Madde 6</b>	Kullanılan tüm dokümanlar güncel ve yeterliydi.	{1,2,3,4,5}
<b>Madde 7</b>	Ders için yeterli sınıf, laboratuvar ve tartışma alanı hazırlanmıştı.	{1,2,3,4,5}
<b>Madde 8</b>	Quiz, ödev, proje ve örnekler eğitime katkısı sağladı.	{1,2,3,4,5}

<b>Madde 9</b>	Derslere katılmaktan çok zevk aldım ve hevesliydim.	{1,2,3,4,5}
<b>Madde 10</b>	Dönem başındaki beklentilerim, dönem sonunda karşılandı.	{1,2,3,4,5}
<b>Madde 11</b>	Kurs kişisel gelişimim için faydalı oldu.	{1,2,3,4,5}
<b>Madde 12</b>	Kurs, hayata yeni bir perspektifte bakmamı sağladı.	{1,2,3,4,5}
<b>Madde 13</b>	Eğitmenin bilgisi yeterli ve günceldi.	{1,2,3,4,5}
<b>Madde 14</b>	Eğitmen sınıfa hazır geldi.	{1,2,3,4,5}
<b>Madde 15</b>	Eğitmen, anlatılan ders planına uydu.	{1,2,3,4,5}
<b>Madde 16</b>	Eğitmen, dersini adanmış ve anlaşılır bir şekilde sundu.	{1,2,3,4,5}
<b>Madde 17</b>	Eğitmen, derslere zamanında geldi.	{1,2,3,4,5}
<b>Madde 18</b>	Eğitmenin konuşması akıcı ve düzgündü.	{1,2,3,4,5}
<b>Madde 19</b>	Eğitmen, ders saatlerini efektif kullandı.	{1,2,3,4,5}
<b>Madde 20</b>	Eğitmen, dersi öğrencilere sevdirmeye konusunda istekliydi.	{1,2,3,4,5}
<b>Madde 21</b>	Eğitmen, öğrencilere karşı pozitif bir yaklaşım sergiliyordu.	{1,2,3,4,5}
<b>Madde 22</b>	Eğitmen, öğrencilerin yorumlarına karşı açık ve saygılıydı.	{1,2,3,4,5}
<b>Madde 23</b>	Eğitmen, öğrencileri derse katılma konusunda cesaretlendirirdi.	{1,2,3,4,5}
<b>Madde 24</b>	Eğitmen, öğrencilere derslerine katkı sağlayacak ödev/projeler verdi.	{1,2,3,4,5}
<b>Madde 25</b>	Eğitmen, öğrencilerin kurs ile ilgili sorularını her ortamda cevapladı.	{1,2,3,4,5}
<b>Madde 26</b>	Eğitmenin kursu kıymetlendirme adına yaptığı tüm sınavlar, kursun amacına yönelikti.	{1,2,3,4,5}
<b>Madde 27</b>	Eğitmen, sınavların çözümlerini öğrencilerle paylaştı ve tartıştı.	{1,2,3,4,5}
<b>Madde 28</b>	Eğitmen, bütün öğrencilere eşit ve tarafsız davrandı.	{1,2,3,4,5}

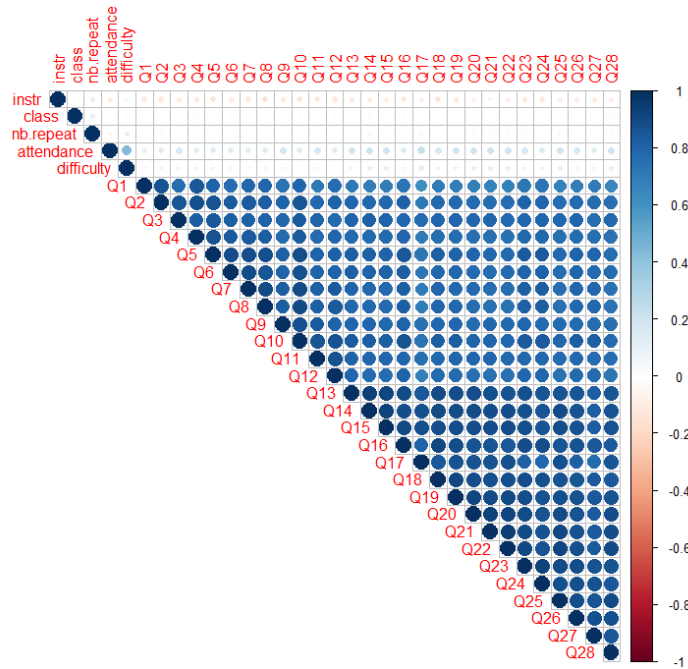
*Madde1-Madde28 aralığındaki tüm öznitelikler Likert tipinde olup {1,2,3,4,5} değerlerinden birini almaktadır.*

Suh (2016) tarafından yayımlanan öğrenme analitiği ve eğitimde veri madenciliği adlı çalışmada aynı veri seti kullanılmış ve “tekrar (nb.repeat)” değişkenini sınıflandırıcı olarak belirlenmiştir. Çalışmada, öğrencinin dersi kaç kez aldığı üzerine bir araştırma yapılmıştır. Sonuç olarak, çoğu öğrencinin dersi “yalnızca bir kez aldığı” sonucuna ulaşmıştır. Bu sonuç

makul bir sonuç olarak değerlendirilmiştir. Ancak dağılımın çok çarpık olması modelin yorumlanabilir, kabul edilebilir bir sınıflandırıcı oluşturma işlemini biraz zorlaştırmıştır. Bu nedenle bahsi geçen yazar tarafından sınıf olarak “eğitmen (instr)” gibi farklı bir sınıf seçilmesi ve bu sınıflar ile diğer nitelikler arasındaki ilişkinin incelenmesi önerilmektedir. Bu doğrultuda bu tez çalışmasında, “eğitmen” değişkeni sınıf (bağımlı değişken) olarak seçilmiş ve üç öğretim elemanı “1”, “2” ve “3” şeklinde sınıflandırılarak veri seti oluşturulmuştur.

### 3.3. Veri Analizi

Makine öğrenmesi yöntemleri geleneksel yöntemlerden farklıdır. MÖ yöntemleri için doğrusallık, homojenlik ya da normallik varsayımlarına ihtiyaç duyulmamaktadır. Makine öğrenmesi, veri madenciliği gibi yöntemler verileri anlamlı bir hale dönüştürmektedir. Bunun için kümeleme, sınıflandırma, tahmin yöntemlerini kullanmaktadır (Filiz, Aşkin ve Öz, 2018; Filiz, Karadağ ve Aşkin, 2018). Yine de makine öğrenmesi çalışmalarında veri setine, veri ön işleme süreci uygulanarak; dağılım, homojenlik, korelasyonel ilişki gibi araştırmalar yapılması önerilmektedir (Shalev-Shwartz ve Ben-David, 2014). Bu aşamada *summary()* ya da daha hızlı bilgi edinmek için *profiling\_num()* fonksiyonu ile verinin betimsel istatistik sonuçları incelenmektedir ve bu çalışma betimsel istatistik sonuçları uygun bulunmuştur. Veri setinde kayıp veri olup olmadığı kontrol etmek için *sum(is.na(veri))* fonksiyonu kullanılmaktadır ve yine bu çalışma için herhangi bir kayıp veri olmadığı bilgisine ulaşılmıştır. Çalışma için korelasyonel ilişki incelenmiş ve bunun için *cor.test()* fonksiyonu kullanılmıştır (Kodların tamamı Ek-1’de yer almaktadır).



Şekil 3. 1. TÖD Veri Seti Korelasyonel İlişkisi

Değişkenler arasındaki korelasyon Şekil 3.1’de yer almaktadır. Burada mavi renkli daireler pozitif korelasyonu, kırmızı renkli daireler ise negatif korelasyonu belirtmektedir. Renklerin koyuluğu korelasyonun derecesini, dairelerin büyüklüğü ise korelasyonun istatistikî önemini ifade etmektedir. Değişkenlerin ilişki oranları ve düzeylerine Tablo 3.2’de yer verilmiştir.

Tablo 3. 2. TÖD Korelasyonel İlişki Sonuçları

TÖD Korelasyonel İlişki								
Değişkenler	Oran	İlişki	Değişkenler	Oran	İlişki	Değişkenler	Oran	İlişki
<b>egitmen</b>	<b>%3.98</b>	<b>NDD</b>	<b>Madde 1</b>	%5.21	PDD	<b>Madde 15</b>	%92.98	PYD
<b>tekrar</b>	%9.15	PDD	<b>Madde 2</b>	%86.61	PYD	<b>Madde 16</b>	%89.48	PYD
<b>katılım</b>	%7.80	NDD	<b>Madde 3</b>	%85.07	PYD	<b>Madde 17</b>	%80.26	PYD
<b>zorluk</b>	%43.67	POD	<b>Madde 4</b>	%82.63	PYD	<b>Madde 18</b>	%84.57	PYD
			<b>Madde 5</b>	%86.78	PYD	<b>Madde 19</b>	%90.45	PYD
			<b>Madde 6</b>	%88.28	PYD	<b>Madde 20</b>	%91.17	PYD
			<b>Madde 7</b>	%89.26	PYD	<b>Madde 21</b>	%92.76	PYD
			<b>Madde 8</b>	%89.99	PYD	<b>Madde 22</b>	<b>%94.14</b>	<b>PYD</b>
			<b>Madde 9</b>	%82.88	PYD	<b>Madde 23</b>	%90.07	PYD
			<b>Madde 10</b>	%87.18	PYD	<b>Madde 24</b>	%92.29	PYD
			<b>Madde 11</b>	%85.65	PYD	<b>Madde 25</b>	%87.66	PYD
			<b>Madde 12</b>	%86.33	PYD	<b>Madde 26</b>	%88.58	PYD
			<b>Madde 13</b>	%79.25	PYD	<b>Madde 27</b>	%87.76	PYD
			<b>Madde 14</b>	%93.58	PYD	<b>Madde 28</b>	%84.65	PYD

(NDD: Negatif Düşük Düzey, POD: Pozitif Orta Düzey, PDD: Pozitif Düşük Düzey, PYD: Pozitif Yüksek Düzey)

Tablo 3.2 incelendiğinde %94.14’lük oran ile *Madde22* değişkeni yüksek korelasyon, %3.98’lik oran ile *sınıf* değişkeni düşük korelasyon göstermektedir.

Çalışmada kullanılmak üzere, Türkiye Öğrenci Değerlendirmesi veri setinde makine öğrenmesi sınıflandırma algoritmalarından; C5.0 Karar Ağacı, Rastgele Orman ve Yapay Sinir Ağları algoritmaları kullanılarak algoritmaların performansları araştırılmıştır. Çalışmada R dilinde kodlar yazılmış ve yine bu dilde yazılmış hazır paketler kullanılmıştır. Kodların gerçekleştirilmesinde geliştirme aracı olarak RStudio ortamından yararlanılmıştır.

### 3.2.1. Karar Ağaçları Uygulaması

Karar ağaçları; CART, CHAID, ID3, C4.5 ve C5.0 gibi algoritmalara sahiptir. Fakat her algoritmanın kendine has belirli özellikleri bulunmaktadır. Bu özellikler incelenerek veri seti için doğru algoritma seçilmektedir. Çalışmada kullanılan veri seti için uygun algoritmanın C5.0 olduğuna karar verilmiştir. Sınıflandırma işleminin gerçekleştirilmesi için “C50” paketinin yüklenmesi ve kütüphaneye erişimin sağlanması gerekmektedir (Kuhn ve diğerleri,2018). Paket *C5.0()* fonksiyonu ile çalıştırılmakta ve *plot()* fonksiyonu ile ağaç modeline erişilmektedir. Karar ağaçlarında yaprak düğümleri ile gözlem sayısının ve oranının doğru sınıflandırılıp sınıflandırılmadığı öğrenilmektedir. Ağacın dallanma aşamasında nasıl bir kural dizisi izlediğini görmek için *summary* fonksiyonu kullanılmaktadır. Summary fonksiyonu verinin yapısını daha iyi anlamak amacıyla kullanılan özet fonksiyondur. Eğitim verisi ile modelin oluşturulmasının ardından *predict()* fonksiyonu kullanılarak test verisi ile tahmin gerçekleştirilmektedir. Tüm işlemlerin ardından *confusionMatrix()* fonksiyonu ile modelin performansı tespit edilmektedir (Kuhn,2017)

### 3.2.2. Rastgele Orman Uygulaması

Modelin oluşturulması için *randomForest()* kütüphanesi kullanılmıştır. Burada karşımıza iki temel değer çıkmaktadır. Bunlar ağaç sayısının (*ntree*) ve özellik (*mtry*) sayısının belirlenmesidir. Burada *ntree* parametresi büyüdükçe (ağaç sayısı arttırıldıkça) tasarlanan modelin hata oranı azalmaktadır (Bulut, 2020). Özellik sayısı belirlenirken bağımsız değişken sayısı dikkate alınmaktadır. Rastgele belirlenen özellik değeri ile OBB hatası hesaplanmaktadır. Hatanın olabilecek en küçük değere sahip olması için başlangıçta rastgele belirlenmiş olan özellik değerinde düzenleme yapılmaktadır. Düzenleme çalışması ile elde edilen sonuçlar incelenmekte ve en düşük hata değerine sahip özellik değeri çalışmada kullanılmaktadır. Ardından eğitim seti için model oluşturulmakta ve oluşturulan modelin kontrol edilmesi için test seti kullanılmaktadır. Elde edilen sonuç ile algoritmanın performansı tespit edilmektedir.



### 3.2.3. Yapay Sinir Ağları Uygulaması

Yapay sinir ağları uygulamasında veri setinin *numeric* formatında olması gerekmektedir. Değişkenlerin her biri standardize edilerek değerler 0 ile 1 arasına sıkıştırılmıştır. Yapay Sinir Ağları için *neuralnet()* kütüphanesi kullanılmakta ve veri seti %70 eğitim ve %30 test (Torgo, 2011) verisi olarak ayrılmaktadır. Bundan sonraki süreçte ağırlık hesaplaması yapılmakta ve bunun için iki yol izlenebilmektedir. İlki gerçek sınıflandırma oranlarının tespit edilmesi ve *hidden()* değerinin değiştirilerek gerçek orana eşit ya da yakın değere ulaşmaya çalışılmasıdır. Burada uygun ağırlık değerinin bulunması için YSA'nın matematiksel formülü koda dönüştürülerek kullanılmaktadır.

```
girdi <- (bias) + (ağırlık *değişken1) + (ağırlık *değişken2)  
girdi  
cıktı <- 1/(1+exp(-girdi))  
cıktı
```

İkinci yol ise, ilgili kodun yazılarak gerçek değer ile tahmin değeri arasındaki minimum farka ulaşana kadar ağırlığın araştırılmasıdır. Bu çalışmanın veri setinde 33 değişken bulunmasından dolayı analiz sürecini kolaylaştıracağı ve zamandan tasarruf sağlayacağı düşünüldüğünden ikinci yol kullanılmıştır. Uygun gizli katman sayısının ayarlanmasının ardından tahmin aşamasına geçilerek hata ve doğruluk değerlerine erişilmektedir.

## BÖLÜM IV

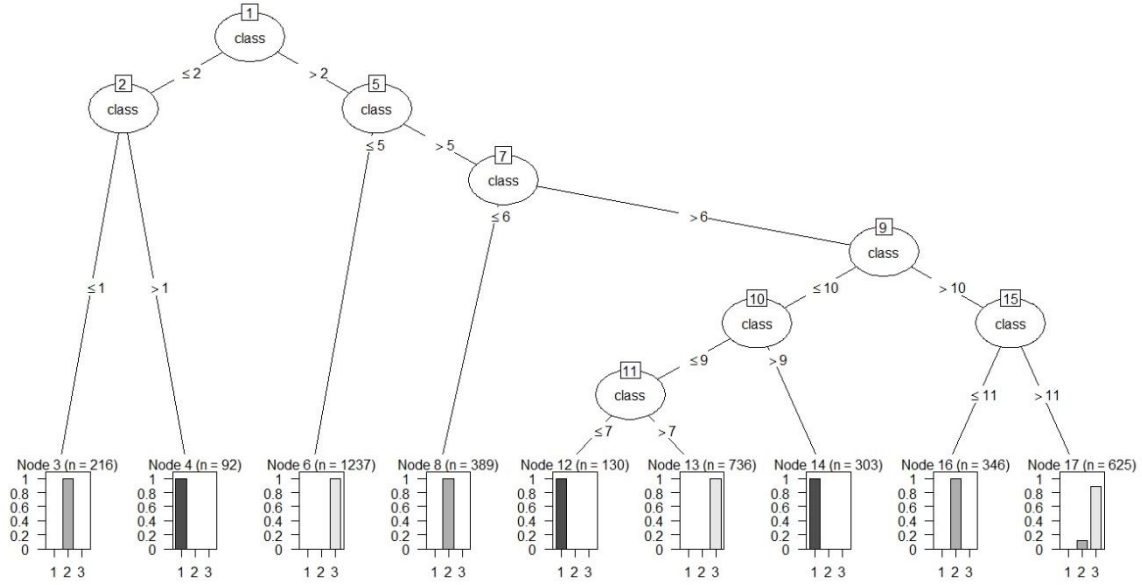
### BULGULAR

Bu bölümde Türkiye Öğrenci Değerlendirme verileri üzerinde gerçekleştirilen C5.0 Karar Ağacı, Rastgele Orman ve Yapay Sinir Ağları algoritmalarına ait elde edilen bulgular ve alt problemlere aranan cevaplara ayrı ayrı yer verilmiştir.

#### 4.1. Algoritmaların Performanslarının İncelenmesi

##### 4.1.1. Karar Ağacı

Karar Ağacı algoritması ile ilgili kodlar Ek-2’de yer almaktadır.



Şekil 4. 1. C5.0 Karar Ağacı Modeli

Eğitim verisi kullanılarak oluşturulan model Şekil 4.1’de yer almaktadır. Şekil incelendiğinde, belirli bir kurala göre dağılımın gerçekleştiği görülmektedir. Node 17’ye odaklanıldığında, diğer yapraklardan farklı bir karar yaprağı gözlemlenmektedir. Bu farklılık dallanmanın Node 17 aşamasında bir hata olduğunu göstermektedir.

Tablo 4.1’de dallanmanın nasıl bir kural izlediğine yer verilmiştir.

Tablo 4. 1. *Karar Ağacı Karar Kuralları*

Kural 1: Eğer sınıf > 9 sınıf ≤ 10 ise “1 Numaralı Eğitimci”	Kural 5: Eğer sınıf > 10 sınıf ≤ 11 ise “2 Numaralı Eğitimci”
Kural 2: Eğer sınıf > 6 sınıf ≤ 7 ise “1 Numaralı Eğitimci”	Kural 6: Eğer sınıf ≤ 1 ise “2 Numaralı Eğitimci”
Kural 3: Eğer sınıf > 1 sınıf ≤ 2 ise “1 Numaralı Eğitimci”	Kural 7: Eğer sınıf > 7 sınıf ≤ 9 ise “3 Numaralı Eğitimci”
Kural 4: Eğer sınıf > 5 sınıf ≤ 6 ise “2 Numaralı Eğitimci”	Kural 8: Eğer sınıf > 2 sınıf ≤ 5 ise “3 Numaralı Eğitimci”
	Kural 9: Eğer sınıf > 11 ise “3 Numaralı Eğitimci”

Tablo 4.1’de C5.0 Karar Ağacı algoritmasının karar kuralları yer almaktadır. Her bir derse özel ismi temsil eden *sınıf* değişkeni en belirleyici nitelik olarak ortaya çıkmıştır. Ayrıca karar kuralları incelendiğinde en önemli değişkenin de *sınıf* değişkeni olduğu görülmektedir.

Karar kurallarının yorumu şu şekilde yapılmaktadır: Eğer ders isminin kodu 9’dan büyük ya da 10’a eşit veya 10’dan küçük ise “1 Numaralı Eğitimci”, eğer ders isminin kodu 6’dan büyük ya da 7’ye eşit veya 7’den küçükse ise “1 Numaralı Eğitimci”, eğer ders isminin kodu 1’den büyük ya da 2’ye eşit veya 2’den küçük ise “1 Numaralı Eğitimci”, eğer ders isminin kodu 5’den büyük ya da 6’ya eşit veya 6’dan küçük ise “2 Numaralı Eğitimci”, eğer ders isminin kodu 10’dan büyük ya da 11’e eşit veya 11’den küçük ise “2 Numaralı Eğitimci”, eğer ders isminin kodu 1’e eşit veya 1’den küçük ise “2 Numaralı Eğitimci”, eğer ders isminin kodu 7’den büyük ya da 9’a eşit veya 9’dan küçük ise “3 Numaralı Eğitimci”, eğer ders isminin kodu 2’den büyük ya da 5’e eşit veya 5’den küçük ise “3 Numaralı Eğitimci”, eğer ders isminin kodu 11’den büyük ise “3 Numaralı Eğitimci” (Çınar, 2019).

C5.0 algoritması için performans değerinin incelenmesi amacıyla eğitim ve test verisi karışıklık matrisinde incelenmiştir. Tablo 4.2’de model için ve Tablo 4.3’te performans değerlendirmesi için elde edilen sonuçlar yer almaktadır.

Tablo 4. 2. Eğitim Seti Hata Tablosu

		GERÇEK		
		1	2	3
TAHMİN	1	525	0	0
	2	0	951	74
	3	0	0	2524

Tablo 4.2 incelendiğinde, 4074 gözlemden oluşan eğitim setinin 74 tanesi hariç diğerlerinin doğru sınıflandırıldığı görülmektedir. %98.18’lik doğruluk değerine sahip eğitim setine bakılarak gerçek sınıflandırma performans değerine ulaşılamaz. Algoritmanın performansını belirlemek için test seti ile aynı işlemler gerçekleştirilir.

Tablo 4. 3. Test Seti Hata Tablosu

		GERÇEK			
		1	2	3	Toplam
TAHMİN	1	250	0	0	250
	2	0	394	25	419
	3	0	0	1077	1077
Toplam		250	394	1102	N=1746

Tablo 4.3 incelendiğinde, 1746 gözlemden oluşan test setinde sınıflandırmanın doğruluk (Accuracy) değeri %98.57, %95 aralıkla güven aralığı (95% CI) 0.9789-0.9907, pozitif sınıfın toplam tahmin oranı (No Information Rate) 0.6312 ve bu model için “3” numaralı öğretim elemanı pozitif ele alındı, anlamlılık değeri(P-Value)  $2,2 \cdot 10^{-16}$ , kapa değeri 0.9733, kesinlik (Precision) 1.00, duyarlılık (Sensitivity) 0.9773, özgüllük (Specificity) 1.00 ve F-ölçütü 0.988 olarak bulunmuştur.

#### 4.1.2. Rastgele Orman

Rastgele Orman algoritması ile ilgili kodlar Ek-3’de yer almaktadır.

Tablo 4. 4. *OBB Hata Tablosu*

<b>GERÇEK</b>			
	<b>1</b>	<b>2</b>	<b>3</b>
<b>1</b>	543	0	0
<b>2</b>	0	952	65
<b>3</b>	0	16	2519

Rastgele Orman algoritması için OBB hata sonuçları Tablo 4.4’te yer almaktadır. Tablo sonuçlarına göre %1.98’lik bir hata değeri bulunmuştur. Özellik (mtry) düzenlenerek hata değeri değiştirilmiş ve Tablo 4.5’te yer alan düzenlenmiş özellik değeri ile veri setinde %1.93’lük OBB hata değeri hesaplanmıştır.

Tablo 4. 5. *Mtry Düzenlenmiş OBB Hata Tablosu*

<b>GERÇEK</b>				
	<b>1</b>	<b>2</b>	<b>3</b>	<b>Sınıfların Hata Değeri</b>
<b>1</b>	543	0	0	0.00
<b>2</b>	0	952	65	0.0639
<b>3</b>	0	14	2521	0.0055

Sınıfların hata oranları tek tek incelendiğinde ise; “1” numaralı öğretim elemanının 0.00, “2” numaralı öğretim elemanının 0.0639 ve “3” numaralı öğretim elemanının 0.0055 şeklinde bulunmuştur.

Tablo 4. 6. Eğitim Seti Hata Tablosu

		GERÇEK		
		1	2	3
TAHMİN	1	543	0	0
	2	0	998	2
	3	0	19	2533

Eğitim seti ile oluşturulan modelin Tablo 4.6’da yer alan sonuçları incelendiğinde, %99.49’luk doğruluk değerine sahip olduğu bulunmuştur. Aynı işlemler test seti için gerçekleştirilerek algoritmanın performansı araştırılmıştır.

Tablo 4. 7. Test Seti Hata Tablosu

		GERÇEK			
		1	2	3	Toplam
TAHMİN	1	232	0	0	232
	2	0	398	5	403
	3	0	29	1061	1090
Toplam		232	427	1066	N=1725

Tablo 4.7 incelendiğinde, 1725 gözlemden oluşan test setinde sınıflandırmanın doğruluk (Accuracy) değeri %98.03, %95 aralıkla güven aralığı (95% CI) 0.9726-0.9863, pozitif sınıfın toplam tahmin oranı (No Information Rate) 0.618 ve bu model için “3” numaralı öğretim elemanı pozitif ele alındı, anlamlılık değeri(P-Value)  $2,2 \cdot 10^{-16}$ , kapa değeri 0.9631, kesinlik (Precision) 0.9734, duyarlılık (Sensitivity) 0.9953, özgüllük (Specificity) 0.9560 ve F-ölçütü 0.9842 olarak bulunmuştur.



Tablo 4. 8. *Yapay Sinir Ağı Model Detayları*

Girdi Katmanı	Bağımsız Değişkenler	sınıf	Madde 7	Madde18	
		tekrar	Madde 8	Madde19	
		katılım	Madde 9	Madde20	
		zorluk	Madde10	Madde21	
		Madde 1	Madde11	Madde22	
		Madde 2	Madde12	Madde23	
		Madde 3	Madde13	Madde24	
		Madde 4	Madde14	Madde25	
		Madde 5	Madde15	Madde26	
		Madde 6	Madde16	Madde27	
			Madde17	Madde28	
			Ara Katman Sayısı	1	
		Ara Katman	Ara Katman İçindeki Bölüm Sayısı	5	
	Aktivasyon Fonksiyonu	Sigmoid Fonksiyonu			
	Bağımlı Değişken	egitmen			
Çıktı Katmanı	Çıkış Katmanı Birim Sayısı	3			
	Aktivasyon Fonksiyonu	Sigmoid Fonksiyonu			

Eğitim seti kullanılarak tasarlanan modelin girdi katmanındaki 32 değişkenin algoritma aracılığı ile ilişkilendirilmesinin ardından tek katmandaki 5 nöron ile gerçek değer ile tahmin değerinin arasında 0.0744'lük bir fark bulunarak, mümkün olan en ideal model tasarlanmıştır. Tasarım aşamasında aktivasyon fonksiyonu olarak Sigmoid Fonksiyonu kullanılmıştır. *egitmen(instr)* bağımlı değişkeni üç farklı öğretim elemanını sınıflandırdığından çıktı biriminde 3 nöron bulunmaktadır.



Tablo 4. 9. *Test Seti Hata Tablosu*

		GERÇEK			
		1	2	3	Toplam
TAHMİN	1	131	53	51	235
	2	2	222	195	419
	3	0	21	1071	1092
Toplam		133	296	1317	N=1746

Tablo 4.9 incelendiğinde, 1746 gözlemden oluşan test setinde sınıflandırmanın doğruluk (Accuracy) değeri %81.55, pozitif sınıfın toplam tahmin oranı (No Information Rate) 0.7542 ve bu model için “3” numaralı öğretim elemanı pozitif ele alındı, anlamlılık değeri(P-Value)  $2,2 \cdot 10^{-16}$ , kesinlik (Precision) 0.9807, duyarlılık (Sensitivity) 0.8132, özgüllük (Specificity) 0.9082 ve F-ölçütü 0.8905 olarak bulunmuştur.

#### 4.2.Algoritmaların Sınıflandırma Performanslarının Karşılaştırılması

Tablo 4.10 incelendiğinde algoritmaların doğruluk oranları; Karar Ağacı %98.57, Rastgele Orman %98.03 ve Yapay Sinir Ağları %81.55 olarak bulunmuştur.

Tablo 4. 10. *Doğruluk Tablosu*

<b>Makine Öğrenmesi Algoritmaları</b>	<b>Tahmin Doğruluğu Performansı (%)</b>
Karar Ağacı	98.57
Rastgele Orman	98.03
Yapay Sinir Ağları	81.55

Tablo 4.11 incelendiğinde, algoritmaların özgüllük oranları; Karar Ağacı %100, Rastgele Orman %95.60 ve Yapay Sinir Ağları %90.82 olarak bulunmuştur.

Tablo 4. 11. *Özgüllük Tablosu*

<b>Makine Öğrenmesi Algoritmaları</b>	<b>Özgüllük Performansı (%)</b>
Karar Ağacı	100
Rastgele Orman	95.60
Yapay Sinir Ağları	90.82

Tablo 4.12 incelendiğinde, algoritmaların duyarlılık oranları; Karar Ağacı %97.73, Rastgele Orman %99.53 ve Yapay Sinir Ağları %98.07 olarak bulunmuştur.

Tablo 4. 12. *Duyarlılık Tablosu*

<b>Makine Öğrenmesi Algoritmaları</b>	<b>Duyarlılık Performansı (%)</b>
Karar Ağacı	97.73
Rastgele Orman	99.53
Yapay Sinir Ağları	98.07

Tablo 4.13 incelendiğinde, algoritmaların kesinlik oranları; Karar Ağacı %100, Rastgele Orman %97.34 ve Yapay Sinir Ağları %81.32 olarak bulunmuştur.

Tablo 4. 13. *Kesinlik Tablosu*

<b>Makine Öğrenmesi Algoritmaları</b>	<b>Kesinlik Performansı (%)</b>
Karar Ağacı	100
Rastgele Orman	97.34
Yapay Sinir Ağları	81.32

Tablo 4.14 incelendiğinde, algoritmaların F-Ölçütü oranları; Karar Ağacı %98.8, Rastgele Orman %98.42 ve Yapay Sinir Ağları %89.05 olarak bulunmuştur.

Tablo 4. 14. *F-Ölçütü Tablosu*

<b>Makine Öğrenmesi Algoritmaları</b>	<b>F-Ölçütü Performansı (%)</b>
Karar Ağacı	98.8
Rastgele Orman	98.42
Yapay Sinir Ağları	89.05

### 4.3.Algoritmaların En Önemli Yordayıcılarının Belirlenmesi ve Manidarlıklarının Karşılaştırılması

Tablo 4.15 incelendiğinde, Karar Ağacı algoritması için *sınıf* değişkeni %100 oran ile en önemli yordayıcı olarak bulunmuştur.

Tablo 4. 15. Karar Ağacı İçin Değişkenlerin Önem Tablosu

Önem Oranı					
Değişkenler	%	Değişkenler	%	Değişkenler	%
tekrar	0	Madde 1	0	Madde 15	0
katılım	0	Madde 2	0	Madde 16	0
zorluk	0	Madde 3	0	Madde 17	0
		Madde 4	0	Madde 18	0
		Madde 5	0	Madde 19	0
		Madde 6	0	Madde 20	0
		Madde 7	0	Madde 21	0
		Madde 8	0	Madde 22	0
		Madde 9	0	Madde 23	0
		Madde 10	0	Madde 24	0
		Madde 11	0	Madde 25	0
		Madde 12	0	Madde 26	0
		Madde 13	0	Madde 27	0
		Madde 14	0	Madde 28	0

Tablo 4.16 incelendiğinde, Rastgele Orman algoritması için *sınıf* değişkeni en önemli yordayıcı olarak bulunmuştur.

Tablo 4. 16. *Rastgele Orman İçin Değişkenlerin Önem Tablosu*

<b>Önem Oranı</b>								
<b>Değişkenler</b>	<b>MDA</b>	<b>MDG</b>	<b>Değişkenler</b>	<b>MDA</b>	<b>MDG</b>	<b>Değişkenler</b>	<b>MDA</b>	<b>MDG</b>
<b>tekrar</b>	95.365	4.404	<b>Madde 1</b>	181.246	5.010	<b>Madde 15</b>	68.720	1.621
<b>katılım</b>	111.392	10.590	<b>Madde 2</b>	119.061	3.026	<b>Madde 16</b>	74.742	1.180
<b>zorluk</b>	-0.6005	9.100	<b>Madde 3</b>	159.661	4.435	<b>Madde 17</b>	128.846	1.818
			<b>Madde 4</b>	93.235	2.292	<b>Madde 18</b>	78.790	2.270
			<b>Madde 5</b>	55.759	1.566	<b>Madde 19</b>	47.966	1.417
			<b>Madde 6</b>	75.952	2.704	<b>Madde 20</b>	78.577	1.547
			<b>Madde 7</b>	117.702	2.885	<b>Madde 21</b>	78.208	1.399
			<b>Madde 8</b>	178.539	3.964	<b>Madde 22</b>	95.885	2.863
			<b>Madde 9</b>	141.069	3.041	<b>Madde 23</b>	49.731	1.488
			<b>Madde 10</b>	122.769	2.727	<b>Madde 24</b>	134.139	2.350
			<b>Madde 11</b>	173.828	3.401	<b>Madde 25</b>	65.280	1.232
			<b>Madde 12</b>	157.897	3.971	<b>Madde 26</b>	31.186	1.563
			<b>Madde 13</b>	87.090	2.704	<b>Madde 27</b>	145.157	2.903
			<b>Madde 14</b>	75.967	2.197	<b>Madde 28</b>	8.161	2.079

Tablo 4.17 incelendiğinde, Yapay Sinir Ağı algoritması için *sınıf* değişkeni en önemli yordayıcı olarak bulunmuştur.

Tablo 4. 17. *Yapay Sinir Ağları İçin Değişkenlerin Önem Tablosu*

<b>Önem Oranı</b>					
<b>Değişkenler</b>	<b>%</b>	<b>Değişkenler</b>	<b>%</b>	<b>Değişkenler</b>	<b>%</b>
<b>tekrar</b>	-103486,90	<b>Madde 1</b>	-39188,53	<b>Madde 15</b>	134853,20
<b>katılım</b>	52633,88	<b>Madde 2</b>	-64706,58	<b>Madde 16</b>	28404,83
<b>zorluk</b>	-75134,17	<b>Madde 3</b>	343544,90	<b>Madde 17</b>	173284,10
		<b>Madde 4</b>	-287912,10	<b>Madde 18</b>	331054,50
		<b>Madde 5</b>	145950,80	<b>Madde 19</b>	-189881,50
		<b>Madde 6</b>	-74900,46	<b>Madde 20</b>	-228801,60
		<b>Madde 7</b>	-101197,20	<b>Madde 21</b>	48888,62
		<b>Madde 8</b>	-79816,24	<b>Madde 22</b>	125586,40
		<b>Madde 9</b>	97558,69	<b>Madde 23</b>	355348,90
		<b>Madde 10</b>	-378196,10	<b>Madde 24</b>	72769,76
		<b>Madde 11</b>	321091,30	<b>Madde 25</b>	-649781,10
		<b>Madde 12</b>	-897,05	<b>Madde 26</b>	688,22
		<b>Madde 13</b>	-208950,50	<b>Madde 27</b>	108620,50
		<b>Madde 14</b>	-76750,45	<b>Madde 28</b>	10690,46

Tablo 4.18 incelendiğinde, öğretim elemanı sınıflamasında Karar Ağacı, Rastgele Orman ve Yapay Sinir Ağları algoritmaları için manidarlık düzeyi  $2,2.10^{-16}$  olarak bulunmuştur.

Tablo 4. 18. *Algoritmaların Manidarlık Düzeyinin Karşılaştırılması*

<b>Makine Öğrenmesi Algoritmaları</b>	<b>Manidarlık</b>
Karar Ağacı	$2,2.10^{-16}$
Rastgele Orman	$2,2.10^{-16}$
Yapay Sinir Ağları	$2,2.10^{-16}$

## BÖLÜM V

### TARTIŞMA, SONUÇ VE ÖNERİLER

Araştırmanın bu bölümünde araştırma sonucu elde edilen sonuçlar tartışılmış ve öneriler üzerinde durulmuştur.

#### 5.1.Sonuç ve Tartışma

Bu araştırma kapsamında, eğitim alanında büyük veri çalışmalarına temel oluşturması amacıyla makine öğrenmesi algoritmalarından hangilerinin alanda kullanılabileceği araştırılmıştır. Bunun için Türkiye Öğrenci Değerlendirmesi gerçek veri setinde, öğrenciler tarafından öğretim elemanı performanslarının belirlenmesi ile ilişkisi olduğu düşünülen derse özel 28 soru ve 5 özellikten oluşan faktörler; Karar Ağacı, Rastgele Orman ve Yapay Sinir Ağları analiz yöntemleri çerçevesinde incelenmiş ve bu üç makine öğrenmesi algoritması sınıflandırma performansları açısından karşılaştırılmıştır.

Makine öğrenmesi denetimli öğrenme algoritmaları ile gerçekleştirilen analizler için veri seti, eğitim ve test verisi olarak ikiye bölünerek kullanılmaktadır (Witten, Frank, Hall, Pal, ve Data, 2005). Bu çalışmada Karar Ağacı, Rastgele Orman ve Yapay Sinir Ağları algoritmaları için veri seti, %70 eğitim verisi ve %30 test verisi (Torgo, 2011) olacak şekilde ele alınmıştır. Çalışma, üç algoritmanın performanslarının incelenmesi ile oluşmaktadır.

Çalışmanın sonuçları incelendiğinde, üç sınıflandırıcı arasından görece Karar Ağaçları daha iyi performans gösteren algoritma olarak tespit edilmiştir. Karar ağacına çok yakın performans göstermesi ile Rastgele Orman algoritması diğer iyi performans gösteren algoritma olarak belirlenmiştir. Yapay Sinir Ağları, diğer iki algoritma ile karşılaştırıldığı düşük performanslı olarak bulunmuştur. Aynı şekilde özgüllük oranı açısından Karar Ağacı algoritması, Rastgele Orman ve Yapay Sinir Ağları algoritmalarına göre daha başarılı bulunmuştur. Duyarlılık oranı daha yüksek olan Rastgele Orman algoritması, Karar Ağacı ve Yapay Sinir Ağları algoritmalarına göre daha başarılı olarak tespit edilmiştir. Kesinlik oranı Karar Ağacı algoritmasında, Rastgele Orman ve Yapay Sinir Ağları algoritmalarına göre daha yüksek sonuçlanmıştır. Son olarak, F-Ölçütü için Karar Ağacı algoritmasının Rastgele Orman ve Yapay Sinir Ağları algoritmalarına göre daha yüksek sonuç verdiği gözlenmiştir. Algoritmaların model performansları incelendiğinde duyarlılık hariç diğer parametreler için Karar Ağacı algoritmasının performansı daha başarılı bulunmuştur.

Karar Ağacı, Rastgele Orman ve Yapay Sinir Ağları algoritmaları model oluşturma ve tahmin süreçlerinde farklı algoritma bağlantıları kurmalarından dolayı farklı performans göstermektedir. Karar Ağaçları belirli bir kural doğrultusunda en iyi özellik üzerinden dallanma



kuralı ile ilerliyorken, Rastgele Orman birçok ağaç arasından seçtiği en önemli özellik üzerinden ilerlemektedir. Yapay Sinir Ağları ise, biyolojik sinir ağlarına benzer yapısı ile nöronlar arası bağlantı kurmaktadır (Amasyalı, 2008; Yılmaz, 2014; Aggarwal, 2018).

Literatürde, eğitim alanında makine öğrenmesi çalışmaları kısıtlı sayıda bulunmaktadır. Bunlardan eğitim alanında makine öğrenmesi sınıflandırma algoritmaları için kullanılanları bu başlık altında araştırılmıştır. Luhaybi ve diğerleri (2018) tarafından sınıflama yöntemleri kullanılarak öğrenci başarısızlığının tahmin edildiği çalışmada, Karar Ağacı algoritması en başarılı olarak belirlenmiştir. Orjuela (2019) tarafından yürütülen Kolombiya Üniversitesi'nin Ulusal Üniversitesi'nde de geçerli olan modeller adlı çalışmada Karar Ağaçları algoritmasının performansı yüksek bulunmuştur. Blasi ve Alsuwaiket (2020) tarafından yürütülen bir çalışmada ise, Karar Ağacı ve YSA algoritmaları kullanılarak yüksek öğrenimde öğrenci suistimallerinin analizi araştırılmış ve sonuç olarak iki algoritmanın da yüksek performans gösterdiği tespit edilmiştir. Kartal ve diğerleri (2019) tarafından öğrencilerin öğrenme stillerinin modellendiği bir araştırma yayımlanmıştır. Karar Ağacı en iyi performans gösteren algoritmalar arasında yer almıştır. Suh (2016) tarafından yayımlanan öğrenme analitiği ve eğitimde veri madenciliği adlı çalışmada *Türkiye Öğrenci Değerlendirmesi Veri Seti* ile öğrencilerin dersi kaç kez aldıkları araştırılmıştır. Bunun için Naive Bayes, k-En Yakın Komşu, Lojistik Regresyon, J4.8 Karar Ağacı, JRip, Çok Katmanlı Algılayıcı ve ZeroR algoritmaları kullanılmış ve Karar Ağacı algoritması en iyi performans gösteren algoritma olarak bulunmuştur. Çifçi, Kaleli ve Günal (2018) tarafından yürütülen çalışmada makine öğrenmesi algoritmalarından C4.5 Karar Ağacı, Naive Bayes, Derin Öğrenme ve K- En Yakın Komşu algoritmaları kullanılmıştır. *Türkiye Öğrenci Değerlendirmesi Veri Seti* ile öğretim elemanı performansının araştırıldığı çalışmada, Derin Öğrenme algoritmasının ardından ikinci sırada %94.57 ile Karar Ağacı en iyi performans gösteren algoritma olarak bulunmuştur. Ahmed, Rizaner ve Ulusoy (2016) tarafından öğretim elemanı performansını tahmin etmek için veri madenciliğini yöntemlerinin araştırıldığı bir çalışmada *Türkiye Öğrenci Değerlendirme Veri Seti* için J48 Karar Ağacı, Çok Katmanlı Algı, Naive Bayes ve Sıralı Minimal Optimizasyon algoritmaları kullanılmış ve Karar Ağacı en iyi performans gösteren algoritma olarak bulunmuştur.

Drăgulescu ve diğerleri (2015) tarafından çok sınıflı sınıflandırma probleminde ödev gönderimlerinin tahmin edilmesi için yürütülen çalışmada en iyi performans gösteren algoritmanın Rastgele Orman algoritması olduğunu tespit edilmiştir. Gök (2017) tarafından akademik başarının tahmin edilmesi sürecinde makine öğrenmesi algoritmalarının kullanılması

adlı çalışmada ise yine Rastgele Orman algoritması en iyi performans gösteren algoritma olarak belirlenmiştir. Ortego (2019) tarafından ilköğretim ve ortaöğretim öğrencileri için okunan kitapların metin okunabilirlik derecesine göre sınıflandırılmasına yönelik çalışmada, ilgili parametreler araştırılmış ve sonuç olarak Rastgele Orman algoritması en iyi performans gösteren algoritma olarak belirlenmiştir. Adnan ve diğerleri (2020) derin öğrenme teknikleriyle mobil öğrencilerin performanslarını araştırmış ve Rastgele Orman algoritmasının en iyi performans gösterdiğini ortaya çıkarmıştır. Richard ve Serrurier (2020) tarafından disleksi ve disgrafi tahmini üzerine yürütülen çalışmada Rastgele Orman en başarılı algoritma olarak belirlenmiştir. Selvi (2020) makine öğrenmesi yöntemlerini kullanarak liseye geçiş sınavlarında öğrenci başarısını tahmin etmek için Rastgele Orman algoritmasının en iyi performans gösterdiğini bulmuştur. Gorbani ve Ghousi (2020) tarafından makine öğrenimi tekniklerini kullanarak öğrencilerin performansı tahmin edilmiş ve Rastgele Orman algoritması en başarılı algoritma olarak belirlenmiştir. Demirhan (2018) tarafından otizm spektrum bozukluğunun belirlenmesi için makine öğrenmesi algoritmaları kullanılarak bir çalışma yürütülmüş ve Rastgele Orman algoritması en iyi performans gösteren algoritma olarak belirlenmiştir.

Babić (2017) tarafından öğrencinin akademik motivasyonunu tahmin etmek için makine öğrenmesi yöntemlerinin kullanıldığı çalışmada Sinir Ağları en başarılı performansa sahip algoritma olarak bulunmuştur. Yan ve Au (2019) tarafından makine öğrenimine dayalı çevrimiçi öğrenmedeki davranışların analizleri üzerine yürütülen çalışmada ise yine Yapay Sinir Ağları en başarılı model olarak bulunmuştur. Altabrawee ve diğerleri (2019) tarafından makine öğrenimi tekniklerini kullanarak öğrencilerin performanslarının tahmin edildiği çalışmada Yapay Sinir Ağları en iyi performans gösteren algoritma olarak raporlanmıştır. Musso ve diğerleri (2020) tarafından akademik yörüngelerde temel eğitim sonuçlarının tahmini üzerine yürütülen bir çalışmada Yapay Sinir Ağları mükemmel performans göstermiştir. Zahour (2020) tarafından akademik ve mesleki rehberlik sorularının otomatik sınıflandırılması için makine öğrenimi yöntemlerinin karşılaştırılması çalışmasında Yapay Sinir Ağları algoritması en iyi performans göstermiştir. Achenie ve diğerleri (2020) tarafından küçük çocuklarda otizm taraması için makine öğrenimi stratejisinin kullanıldığı çalışmada Yapay Sinir Ağları yüksek performans gösteren algoritma olarak belirlenmiştir. Yıldız ve Börekci (2020) tarafından makine öğrenmesi yöntemleri ile akademik başarı araştırılmıştır. Yapay Sinir Ağları en iyi performans gösteren algoritma olarak bulunmuştur.

Araştırmalar incelendiğinde, genel olarak Karar Ağacı algoritması, Rastgele Orman ve Yapay Sinir Ağı algoritmalarına göre daha başarılı sonuçlar sergilemiştir. Bu çalışmada ise,

Karar Ağaçları en iyi performans gösteren analiz yöntemi olarak bulunmuştur. Rastgele Orman algoritması ise çok küçük bir farkla iyi performans gösteren diğer algoritmadır. İki algoritmanın da iyi performans gösterme nedeni; Karar Ağacı algoritmasının bir kararın olası tüm sonuçlarını göstermek için dallanma yöntemini kullanması, Rastgele Orman algoritmasının ise tüm karar ağaçlarının çıktılarını dayanarak nihai sonucu veren bir karar ağacı kümesi oluşturmasıdır.

Çalışmada kullanılan veri seti için korelasyon analizi yapılmıştır. Bu doğrultuda; negatif düşük korelasyon gösteren *sınıf*, *tekrar*, *katılım* ve *Madde1* değişkenleri arasında en düşük korelasyonel ilişkiye sahip değişken *sınıf* bağımsız değişkeni olarak bulunmuştur. *zorluk* değişkeni pozitif orta düzey ve *Madde2*, *Madde3*, *Madde4*, *Madde5*, *Madde6*, *Madde7*, *Madde8*, *Madde9*, *Madde10*, *Madde11*, *Madde12*, *Madde13*, *Madde14*, *Madde15*, *Madde16*, *Madde17*, *Madde18*, *Madde19*, *Madde20*, *Madde21*, *Madde22*, *Madde23*, *Madde24*, *Madde25*, *Madde26*, *Madde27*, *Madde28* değişkenleri pozitif yüksek düzeyde ilişki göstermiştir. Yüksek ilişki gösteren değişkenler arasında *Madde28* bağımsız değişkeni en yüksek korelasyonel ilişki gösteren değişken olarak tespit edilmiştir.

Karar Ağacı, Rastgele Orman ve Yapay Sinir Ağı algoritmaları için *sınıf* değişkeni diğer değişkenlere göre en önemli yordayıcı olarak belirlenmiştir. Karar Ağacı ve Rastgele Orman algoritmaları en iyi ve en önemli özellik üzerinden dallanarak ilerlemektedir. Dallanma için en önemli özellik *sınıf* olarak belirlenmiştir. Yapay Sinir Ağları sonuç matrisini etkileyen en önemli değişken ise yine *sınıf* olarak bulunmuştur.

Karar Ağacı, Rastgele Orman ve Yapay Sinir Ağları yöntemlerine göre belirlenen manidarlık, üç algoritma için de aynı düzeydedir ve birbirlerine göre farklılaşmamaktadır.

Bu çalışmanın sonuçları eğitim alanında kullanılabilecek makine öğrenmesi algoritmalarını göstermektedir. Bu çalışmanın sonuçları eğitim alanında kullanılabilecek makine öğrenmesi algoritmalarını göstermektedir. Bunun için Türkiye Öğrenci Değerlendirmesi gerçek veri seti kullanılmıştır. Bahsi geçen veri seti için; Naive Bayes, k-En Yakın Komşu, Lojistik Regresyon, J4.8 Karar Ağacı, JRip, Çok Katmanlı Algılayıcı, ZeroR, C4.5 Karar Ağacı, Derin Öğrenme, Çok Katmanlı Algı ve Sıralı Minimal Optimizasyon algoritmalarının öğretim elemanı sınıflandırma performansları incelenmiştir. Bu doğrultuda çalışmada, C5.0 Karar Ağaçları, Rastgele Orman ve Yapay Sinir Ağlarının öğretim elemanı sınıflandırma performansları incelenmiş ve öğretim elemanı kalitesinin belirlenmesi için bahsi geçen algoritmaların nasıl performans sergilediği araştırılmıştır. Ulusal ve uluslararası

arařtırmalar ile benzerlik gösteren bu alıřma, var olan arařtırmalara bir yenisini ekleyerek literatüre katkı saęlamıřtır.

## 5.2.Öneriler

### 5.2.1. Uygulayıcılar İin Öneriler

- alıřmada ders isminin kodu (sınıf) en önemli deęiřken olarak gözlenmiřtir. Farklı dersler için öęretim elemanlarının sınıflandırıldıęı arařtırmalarda dersin isminin belirlenmesine dikkat edilmelidir.
- Veri setinden deęiřkenin atılması ya da eklenmesi için korelasyon dikkate alınmamalıdır. Bu alıřmada *sınıf* deęiřkeni en düşük korelasyon gösterirken, algoritmalar için en önemli deęiřken olarak tespit edilmiřtir.
- Ulusal ya da uluslararası düzeyde uygulanan sınavların analizinde model yorumlama ařamasını kolaylařtırmak ve zamandan tasarruf saęlamak amacıyla Karar Aęaçları ya da Rastgele Orman algoritmalarının kullanılması önerilmektedir.
- Doğruluk, duyarlılık, özgüllük ve kesinlik gibi deęerler dikkate alındıęında bu veri setine benzer veri setleri ve sınıflandırmalar için Karar Aęacı'nın kullanılması önerilmektedir.
- Rastgele Orman alıřacak uygulayıcılar için uygun mtry deęerinin belirlenmesi sırasında arařtırmaya veri setinin baęımsız deęiřken sayısı kadar deęer vererek arařtırmaya bařlamaları önerilmektedir.

### 5.2.2. Arařtırmacılar İin Öneriler

- Bu alıřma Gazi Üniversitesi öęrencilerinden toplanan, Türkiye Öęrenci Deęerlendirmesi verilerinden oluřmaktadır. PISA, TIMMS ve ABİDE verileri kullanılarak, bu veri setlerine uygun makine öęrenmesi algoritmaları arařtırılabilir.

## KAYNAKÇA

- Achenie, L. E., Scarpa, A., Factor, R. S., Wang, T., Robins, D. L. ve McCrickard, D. S. (2019). A machine learning strategy for autism screening in toddlers. *Journal Of Developmental and Behavioral Pediatrics: JDBP*, 40(5), 369.
- Adnan, M., Habib, A., Ashraf, J., Shah, B. ve Ali, G. (2020). Improving m-learners' performance through deep learning techniques by leveraging features weights. *IEEE Access*, 8, 131088-131106.
- Afrin, F., Rahaman, M. S. ve Hamilton, M. (2020). Mining student responses to infer student satisfaction predictors. arXiv preprint arXiv:2006.07860.
- Aggarwal, C. C. (2018). *Neural networks and deep learning*. Springer, 10, 978-3.
- Ağyar, Z. (2016). Yapay Zekâ ve Sinir Ağları. *Hosting Dergi*.  
<https://www.hostingdergi.com.tr/yapay-zeka-ve-sinir-aglari/> adresinden 14.07.2021'de alınmıştır.
- Ahmed, A. M., Rizaner, A. ve Ulusoy, A. H. (2016). Using data mining to predict instructor performance. *Procedia Computer Science*, 102, 137-142.
- Akpınar, H. (2000). Veri tabanlarında bilgi keşfi ve veri madenciliği. *İ.Ü. İşletme Fakültesi Dergisi*, 29(1), 1-22.
- Al, M., Tucker, A. ve Yousefi, L. (2018). The prediction of student failure using classification methods: a case study. *In Dalam Proc. Int. Conf. Image Process. Pattern Recognit, Hal*, 79-90.
- Alan, A. ve Karabatak, M. (2020). Veri seti-sınıflandırma ilişkisinde performansa etki eden faktörlerin değerlendirilmesi. *Fırat Üniversitesi Mühendislik Bilimleri Dergisi*, 32(2), 531-540.
- AL-Fakhry, N. A. (2016). Summarizing data by using data mining techniques a comparative by using C4. 5 and C5. 0 algorithms. *International Education and Research Journal*, 2(4), 8-12.
- Alpaydin, E. (2020). *Introduction to machine learning* (2.baskı). MIT press.

- Altabrawee, H., Ali, O. A. J. ve Ajmi, S. Q. (2019). Predicting students' performance using machine learning techniques. *Journal of university of babylon for pure and applied sciences*, 27(1), 194-205.
- Amasyalı, M. F. (2008). *Yeni makine öğrenmesi metotları ve ilaç tasarımı uygulamaları*. Doktora Tezi. Yıldız Teknik Üniversitesi, Fen Bilimleri Enstitüsü, İstanbul.
- Asilkan, Ö. ve Irmak, A. G. S. (2009). İkinci el otomobillerin gelecekteki fiyatlarının yapay sinir ağları ile tahmin edilmesi. *Süleyman Demirel Üniversitesi İktisadi ve İdari Bilimler Fakültesi Dergisi*, 14(2), 375-391.
- Ay, Ş. (2020). Model Performansını Değerlendirmek. <https://medium.com/deep-learning-turkiye/model-performans%C4%B1n%C4%B1-de%C4%9Ferlendirmek-metrikler-cb6568705b1> adresinden 14.07.2021'de alınmıştır.
- Babić, I. D. (2017). *Machine learning methods in predicting the student academic motivation*. *Croatian Operational Research Review*, 443-461.
- Başol, G. (2015). *Eğitimde ölçme ve değerlendirme*. Ankara: Pegem Akademi Yayıncılık.
- Bayrakçı, S. ve Albayrak, M. A. (2019). A comparative database survey on the use of big data in academic studies. *Academic Journal of Information Technolog*, 10(36), 73.
- Bell, J. (2020). *Machine learning: hands-on for developers and technical professionals*. John Wiley and Sons. [https://books.google.com.tr/books?hl=tr&lr=&id=p\\_ODwAAQBAJ&oi=fnd&pg=PR27&dq=Bell,+J.+\(2020\).+Machine+learning:+handson+for+developers+and+technical+professionals.+John+Wiley+and+Sons.&ots=mLdthlOaA1&sig=dgOY11NjUVfikP1wP8m3KhqY\\_o&redir\\_esc=y#v=onepage&q=Bel1%2C%20J.%20\(2020\).%20Machine%20learning%3A%20handson%20for%20developers%20and%20technical%20professionals.%20John%20Wiley%20and%20Sons.&f=false](https://books.google.com.tr/books?hl=tr&lr=&id=p_ODwAAQBAJ&oi=fnd&pg=PR27&dq=Bell,+J.+(2020).+Machine+learning:+handson+for+developers+and+technical+professionals.+John+Wiley+and+Sons.&ots=mLdthlOaA1&sig=dgOY11NjUVfikP1wP8m3KhqY_o&redir_esc=y#v=onepage&q=Bel1%2C%20J.%20(2020).%20Machine%20learning%3A%20handson%20for%20developers%20and%20technical%20professionals.%20John%20Wiley%20and%20Sons.&f=false) adresinden 13.07.2021'de alınmıştır.
- Bingöl, K., Aslı, E. R., Örmecioğlu, H. T. ve Arzu, E. R. (2020). Depreme dayanıklı mimari tasarımda yapay zekâ uygulamaları: Derin öğrenme ve görüntü işleme yöntemi ile düzensiz taşıyıcı sistem tespiti. *Gazi Üniversitesi Mühendislik Mimarlık Fakültesi Dergisi*, 35(4), 2197-2210.

- Blasi, A. H. ve Alsuwaiket, M. (2020). Analysis of students' misconducts in higher education using decision tree and ann algorithms. *Engineering, Technology and Applied Science Research, 10*(6), 6510-6514.
- Bontempi, G., Taieb, S. B. ve Le Borgne, Y. A. (2012, July). *Machine learning strategies for time series forecasting*. In European Business Intelligence Summer School, 62-77.
- Boucheron, B. ve Tagliaferri, L. (2019). *Python machine learning projects*. USA: DigitalOcean.
- Breiman, L. (2001). Random forests. *Machine Learning, 45*(1), 5-32.
- Brierley, P., Vogel, D. ve Axelrod, R. (2011). Heritage provider network health prize round 1 milestone prize: How we did it–team ‘market makers’.
- Brownlee, J. (2015). Contract faculty in canada: using access to information requests to uncover hidden academics in canadian universities. *Higher Education, 70*(5), 787-805.
- Brynjolfsson, E., Hitt, L.M. ve Kim, H.H. (2011). Strength in numbers: how does data-driven decisionmaking affect firm performance. *Social Science Research Network, Working Paper Series*.
- Bulut, T. (2020). R programlama diliyle sınıflandırma problemlerinin çözümünde rastgele orman algoritması üzerine bir vaka çalışması: a case study on random forest (rf) algorithm in solving classification problems with r programming language. <https://tevfikbulut.com/2020/05/14/rastgele-orman-algoritmasina-uzerine-bir-vaka-calismasi-a-case-study-on-random-forest-rf-algorithm/> adresinden 14.05.2021’de alınmıştır.
- Büyüköztürk, Ş., Çokluk, Ö. ve Köklü, N. (2019). Sosyal bilimler için istatistik. Ankara: Pegem Akademi.
- Cackett, D. (2013). Information management and big data: a reference architecture. *Oracle: Redwood City, CA, USA*.
- Carvalho, D. R. ve Freitas, A. A. (2004). A hybrid decision tree/genetic algorithm method for data mining. *Information Sciences, 163*(1-3), 13-35.

- Castro, F., Vellido, A., Nebot, A. ve Mugica, F. (2007). Applying data mining techniques to e-learning problems. *In Evolution Of Teaching And Learning Paradigms In Intelligent Environment*, 183-221.
- Chouseinoglou, O. ve Şahin, İ. (2019). Metin madenciliği, makine ve derin öğrenme algoritmaları ile web sayfalarının sınıflandırılması. *Yönetim Bilişim Sistemleri Dergisi*, 5(2), 5-6.
- Cutler, A., Cutler, D. R. ve Stevens, J. R. (2012). *Random forests*. In Ensemble machine learning, 157-175.
- Çakır, Ö. (2008) *Veri madenciliğinde sınıflandırma yöntemlerinin karşılaştırılması: bankacılık müşteri veri tabanı üzerinde bir uygulama*. (Yayımlanmamış Doktora Tezi). Marmara Üniversitesi Sosyal Bilimler Enstitüsü, İstanbul.
- Çınar, A. (2019). Veri madenciliğinde sınıflandırma algoritmalarının performans değerlendirmesi ve r dili ile bir uygulama. *Öneri Dergisi*, 14(51), 90-111.
- Çınar, U. K. (2018). Yapay Sinir Ağları: <https://www.veribilimiokulu.com/blog/yapay-sinir-aglari/> adresinden 14.06.2021'de alınmıştır.
- Çifçi, F., Kaleli, C. ve Günel, S. (2018). Öznitelik seçme ve makine öğrenmesi yöntemleriyle eğitim performansının tahmin edilmesi. *Anadolu Journal of Educational Sciences International*, 8(2), 419-440.
- Das, S., Dey, A., Pal, A. ve Roy, N. (2015). Applications of artificial intelligence in machine learning: review and prospect. *International Journal of Computer Applications*, 115(9).
- Daumé, H. (2017). *A course in machine learning*. Hal Daumé III, 149-155.
- Davenport, T. H. (2014). *Big data at work: dispelling the myths, uncovering the opportunities*. Harvard Business School Publishing Corporation.
- Davenport, T. H., Jeanne G. H. ve Morison. R. (2010). *Analytics at work: smarter decisions, better results*. Harvard Business Press, Massachusetts.
- Dean, J. (2014). *Big data, data mining, and machine learning: value creation for business leaders and practitioners*. John Wiley and Sons.
- Demirhan, A. (2018). Performance of machine learning methods in determining the autism spectrum disorder cases. *Mugla Journal of Science and Technology*, 4(1), 79-84.



- DeVore, S., Yang, J., ve Stewart, J. (2020). Extending machine learning to predict unbalanced physics course outcomes. arXiv preprint arXiv:2002.01964.
- Dhar, V. (2013). *Data science and prediction*. Communications of the ACM, 56(12), 64–73.
- Diebold, F. X. (2003). Big data dynamic factor models for macroeconomic measurement and forecasting. In *Advances in Economics and Econometrics: Theory and Applications, Eighth World Congress of the Econometric Society*, 115-122.
- Dixit, A. (2017). *Ensemble machine learning: a beginner's guide that combines powerful machine learning algorithms to build optimized models*. Packt Publishing Ltd.
- Domingos, P. (2012). A few useful things to know about machine learning. *Communications of the ACM*, 55(10), 78-87.
- Drabkin, R. (2017). Machine learning: the “next big thing” in education.  
<http://www.gettingsmart.com/2017/04/next-big-thing-education/adresinden>  
13.06.2021’de alınmıştır.
- Drăgulescu, B., Bucos, M. ve VasIU, R. (2015). Predicting assignment submissions in a multi-class classification problem. *TEM Journal*, 4(3), 244.
- Ee Publishers (2019). *Techonology and Business for Development*.  
<https://www.ee.co.za/article/application-of-machine-learning-algorithms-in-boiler-plant-root-cause-analysis.html/application-of-machine-learning-algorithms-in-boiler-plant-root-cause-analysis-fig-1> adresinden 13.06.2021’de alınmıştır.
- Efe, E., Bek, Y. ve Şahin, M. (2000). *SPSS’te çözümleri ile istatistik yöntemler II*. Kahramanmaraş Sütçü İmam Üniversitesi Rektörlüğü Yayın No: 73, Ders Kitapları Yayın No: 9, KS Ü. Basımevi, Kahramanmaraş, 214.
- Erl, T., Khattak, W. ve Buhler, P. (2016). *Big data fundamentals: concepts, drivers and techniques*. Boston: Prentice Hall.
- Eynon, R. (2013). The rise of Big Data: what does it mean for education, technology, and media research?
- Freitas, R. D. C. M. (2004). A política de combate à pobreza e as agências multilaterais: um estudo comparativo entre o Brasil e o México nas décadas de 80 e 90.

- Ghorbani, R. ve Ghousi, R. (2020). Comparing different resampling methods in predicting Students' performance using machine learning techniques. *IEEE Access*, 8, 67899-67911.
- Gök, M. (2017). Makine öğrenmesi yöntemleri ile akademik başarının tahmin edilmesi. *Gazi Üniversitesi Fen Bilimleri Dergisi Part C: Tasarım ve Teknoloji*, 5(3), 139-148.
- Gunduz, G. ve Fokoue, E. (2013). UCI machine learning repository [http://mllearn.ics.uci.edu/MLRepository.html]. Irvine, CA: University of California, School of Information and Computer Science.
- Gunduz, N. ve Fokoue, E. (2015). Pattern discovery in students' evaluations of professors: a statistical data mining approach. arXiv preprint arXiv:1501.02263.
- Gümüştaş, E. (2019). *Kayıp gözlem içeren dengesiz veri setlerinin topluluk öğrenme algoritmaları ile sınıflandırılması*. Yüksek Lisans Tezi, Mimar Sinan Güzel Sanatlar Üniversitesi, Fen Bilimleri Enstitüsü, İstanbul.
- Hagan, M. T., Demuth, H. B. ve Beale, M. (1997). *Neural network design*. PWS Publishing Co.
- Hamilton, B. A. (2015). The Field Guide to Data Science, 126. Retrieved from papers3. publication/uuid/1941BECE-325A-45B6-B10C-5A850FA5D609.
- Hao, Q., Galyardt, A., Barnes, B., Branch, R. M. ve Wright, E. (2018). Automatic identification of ineffective online student questions in computing education. In 2018 IEEE Frontiers in Education Conference (FIE) (1-5). IEEE.
- Howell, D.G. (1987). *Statistical methods for psychology (second edition)*. Boston: Duxbury Press.
- Hsieh, W. W. (2009). *Machine learning methods in the environmental sciences: neural networks and kernels*. Cambridge University Press.
- İntellipaat (2016). 10 Big Data Applications in Real Life. <https://intellipaat.com/blog/7-big-data-examples-application-of-big-data-in-real-life/> adresinden 12.05.2021'de alınmıştır.
- James, G., Witten, D., Hastie, T. ve Tibshirani, R. (2013). *An introduction to statistical learning* 112(18). New York: springer.
- Jankowski, N., Duch, W. ve Grąbczewski, K. (2011). *Meta-learning in computational intelligence*. Springer.

- Javed, D. (2019). Big data data mining and machine learning.
- Jordan, M. I. ve Mitchell, T. M. (2015). *Machine learning: trends, perspectives, and prospects*. Science, 255-260.
- Jordan, M., Kleinberg, J. ve Schölkopf, B. (2006). *Information science and statistics*.
- Junior, M. A. D. C. O. (2011). Redes neurais de hopfield para roteamento de redes de comunicação em fpga. *Trabalho de Conclusão de Curso-Engenharia da Computação*, 5.
- Kabathova, J. ve Drlik, M. (2021). Towards predicting student's dropout in university courses using different machine learning techniques. *Applied Sciences*, 11(7).
- Karakuzu, C. (1998). *Yapay sinir ağları ile bir kontrol uygulaması*. Yüksek Lisans Tezi, Kocaeli Üniversitesi Fen Bilimleri Enstitüsü, Kocaeli.
- Kartal, E., Köse Biber, S., Biber, M., Özyaprak, M., Şimşek, İ. ve Can, T. (2019). Makine öğrenmesi tekniklerini ve kolb öğrenme stilleri envanterini kullanarak öğrencilerin öğrenme stillerinin belirlenmesi için bir model önerisi. *Kastamonu Eğitim Dergisi*, 27 (5), 1875-1892.
- Kaya, İ., Oktay, S. ve Engin, O. (2005). Kalite kontrol problemlerinin çözümünde yapay sinir ağlarının kullanımı. *Erciyes Üniversitesi Fen Bilimleri Enstitüsü Dergisi*, 21(1-2), 95.
- Kılıç, S. (2015). Kappa Testi. *Journal of Mood Disorders*, 5(3).
- Kıran, Z. B. (2010). *Lojistik regresyon ve cart analizi teknikleriyle sosyal güvelik kurumu ilaç provizyon sistemi verileri üzerinde bir uygulama*. Yüksek Lisans Tezi, Gazi Üniversitesi, Fen Bilimleri Enstitüsü, Ankara.
- Kondakçı, Y., Emil, S. ve Beycioğlu, K. (2019). *14. uluslararası eğitim yönetimi kongresi tam metin bildiri kitabı*. Eğitim Yöneticileri ve Eğitim Denetçileri Derneği.
- Köktürk, F. (2012). *K-en yakın komşuluk, yapay sinir ağları ve karar ağaçları*. Doktora Tezi, Bülent Ecevit Üniversitesi, Sağlık Bilimleri Enstitüsü, Biyoistatistik Anabilim Dalı. Zonguldak.
- Krichevsky, M., Martynova, J. ve Budagov, A. (2019). Methods of machine learning in the master's educational program. In E3S Web of Conferences (Vol. 135, p. 03069). EDP Sciences.

- Kotsiantis, S. B. (2012). Use of machine learning techniques for educational proposes: a decision support system for forecasting students' grades. *Artificial Intelligence Review*, 37(4), 331-344.
- Kuhn, M. (2017). Classification and Regression Training, Package 'caret' Version 6.0-77. <https://cran.r-project.org/web/packages/caret/caret.pdf> adresinden 10.06.2021'de alınmıştır.
- Kumar, A. (2020). Data Analytics. <https://vitalflux.com/hold-out-method-for-training-machine-learning-model/> adresinden 10.06.2021'de alınmıştır.
- Lakkaraju, H, Aguiar, E, Shan, C, Miller, D, Bhanpuri, N, Ghani, R., vd. (2015). *A machine learning framework to identify students at risk of adverse academic outcomes*. KDD, 1909–1918.
- Latif, S., XianWen, F. ve Wang, L. L. (2021). Intelligent decision support system approach for predicting the performance of students based on three-level machine learning technique. *Journal of Intelligent Systems*, 30(1), 739-749.
- Lee, K., Chung, J. ve Suh, C. (2017). Large-scale and interpretable collaborative filtering for educational data. ML4ED KDD Work, 1-7.
- Liaw, A., ve Wiener, M. (2002). Classification and regression by randomForest. *R news*, 2(3), 18-22.
- Lin, N. (1976). *Foundations of social research*. USA: McGraw-Hill.
- Linn, R. L. ve Gronlund, N. E. (2000). *Measurement and evaluation in teaching*. New York: Macmillan Publishing.
- Luhaybi, M.A., Tucker, A. ve Yousefi, L. (2018). The prediction of student failure using classification methods: a case study. *Computer Science and Information Technology*.
- Maimon, O. Z. ve Rokach, L. (2007). *Data mining with decision trees: theory and applications* (69). World Scientific.
- Maimon, O. ve Rokach, L. (2010). *Data mining and knowledge discovery handbook*. New York, USA.

- Manyika, J., Chui, M., Brown, B., Bughin, J., Dobbs, R., Roxburgh, C., vd. (2011). *Big data: the next frontier for innovation, competition, and productivity*. McKinsey Global Institute.
- Mayer-Schönberger, V. ve Cukier, K. (2013). *Big data: a revolution that will transform how we live, work, and think*. Houghton Mifflin Harcourt.
- McMillan, J. H. (2008). *Assessment essentials for standards-based education*. Corwin Press.
- Mduma, N., Kalegele, K. ve Machuve, D. (2019). A survey of machine learning approaches and techniques for student dropout prediction.
- MEB (2018). PISA. [http://pisa.meb.gov.tr/?page\\_id=18](http://pisa.meb.gov.tr/?page_id=18) adresinden 07.07.2021'de alınmıştır.
- Mohri, M., Rostamizadeh, A. ve Talwalkar, A. (2018). *Foundations of machine learning*. MIT press.
- Molnar, C. (2020). *Interpretable machine learning*. Lulu.com.
- Monino, J. L. ve Sedkaoui, S. (2016). *Big data, open data and data development* (3.baskı).
- Mosquera Navarro, R. (2021). *Sistema de clasificación basado en técnicas inteligentes para identificar el grado de riesgo psicosocial en docentes de educación básica primaria y secundaria en colegios públicos de Colombia*. Doktora Tezi, Kolombiya Üniversitesi.
- Musso, M. F., Rodríguez Hernández, C. F. ve Cascallar, E. C. (2020). Predicting key educational outcomes in academic trajectories: a machine-learning approach.
- Nasiriany, S., Thomas, G., Wang, W., Yang, A., Listgarten, J. ve Sahai, A. (2019). *A comprehensive guide to machine learning*. University of California at Berkeley, 82-88.
- Nasuhoglu, H. (2019). *Eczacılık sektöründe yapay sinir ağları ve zaman serileri analizi ile talep tahmini*. Yüksek Lisans Tezi, T.C. Maltepe Üniversitesi, Fen Bilimleri Enstitüsü, İstanbul.
- Nunn, S, Avella, JT, Kanai, T ve Kebritchi, M. (2016). Learning analytics methods, benefits, and challenges in higher education: a systematic literature review. *Online Learning*, 20(2): 13–29.
- Olson, D. L. ve Delen, D. (2008). *Advanced data mining techniques*. Springer Science and Business Media.

- Ortego, R. G. ve Sánchez, I. M. (2019). Relevant parameters for the classification of reading books depending on the degree of textual readability in primary and compulsory secondary education (cse) students. *IEEE Access*.
- Osmanoğlu, U., Atak, O., Çağlar, K., Kayhan, H. ve Can, T. (2020). Sentiment analysis for distance education course materials: a machine learning approach. *Journal of Educational Technology and Online Learning*, 3(1), 31-48.
- Otálora Orjuela, F. (2019). Modelo para la identificación de patrones de desempeño académico estudiantil para fortalecer el acompañamiento académico en la Universidad Nacional de Colombia.
- Öztemel, E. (2003). *Yapay sinir ağları*. İstanbul: Papatya Yayıncılık.
- Pehlivan, G. (2006). *Chaid analizi ve bir uygulama*. Yüksek Lisans Tezi, Yıldız Teknik Üniversitesi, İstanbul.
- Peng, H. (2013). *Air Quality Prediction by Machine Learning Methods*. The Degree of Master of Science. the University of British Columbia.
- Poole, D. L. ve Mackworth, A. K. (2010). *Artificial Intelligence: foundations of computational agents*. Cambridge University Press.
- Prince, S. J. (2012). *Computer vision: models, learning, and inference*. Cambridge University Press.
- Provost, F. ve Fawcett, T. (2013). *Data science and its relationship to big data and data-driven decision making*. *Big Data*, 1(1), 51-59.
- Qi, Y. (2012). *Random forest for bioinformatics*. In *Ensemble Machine Learning*, 307-323.
- Quinlan, J.R. (1986). Induction of decision trees. *Machine Learning*, 81-106.
- Raj, R. (2021). Enjoy Algorithms. <https://www.enjoyalgorithms.com/blogs/supervised-unsupervised-and-semisupervised-learning> adresinden 10.03.2021'de alınmıştır.
- Rajapakse, J. C. ve Omondi, A. (2006). *FPGA implementations of neural networks*. The Netherlands.
- Richard, G. ve Serrurier, M. (2020). Dyslexia and dysgraphia prediction: a new machine learning approach. arxiv preprint.

- Schapiro, R. E. ve Freund, Y. (2012). Foundations of machine learning.
- Scism, L. ve Maremont, M. (2010). Insurers test data profiles to identify risky clients. *Wall Street Journal*.  
<https://www.wsj.com/articles/SB10001424052748704648604575620750998072986>  
adresinden 13.06.2021'de alınmıştır.
- Selvi, A. (2020). *Bilecik ilinde ilköğretimden liseye geçiş sınavlarında makine öğrenmesi yöntemleri ile öğrenci başarısının tahmini*. Yüksek Lisans Tezi, Şeyh Edebali Üniversitesi, Fen Bilimleri Enstitüsü, Bilecik.
- Shalev-Shwartz, S. ve Ben-David, S. (2014). *Understanding machine learning: from theory to algorithms*. USA: Cambridge University Press, 449.
- Shishodia, K., Sekhon, G. ve Rajpal, P. (2006). An artificial neural network for modeling reliability, availability and maintainability of a repairable system. *Reliability Engineering and System Safety*, (91), 811.
- Siemens, G. ve Long, P. (2011). Penetrating the fog: analytics in learning and education. *Educause review*, 46(5), 30.
- Sitorus, L. (2015). *Algoritma dan pemrograman*. Penerbit Andi.
- Smola, A. ve Vishwanathan, S. (2008). *Introduction to machine learning*. Press Syndicate Of The University Of Cambridge The Pitt Building, Unites Kingdom, Cambridge, 226.
- Sobnath, D., Kaduk, T., Rehman, I. U. ve Isiaq, O. (2020). Feature selection for UK disabled students' engagement post higher education: a machine learning approach for a predictive employment model.
- Solanki, K. ve Dhankar, A. (2017). A review on machine learning techniques. *International Journal of Advanced Research in Computer Science*, 8(3).
- Suh, S. (2016). Learning analytics and educational data mining.
- Şahin, T. (2019). Yapay zekâ yolculuğunda sorular ve cevaplar.  
<https://medium.com/baybaynakit/yapay-zeka-yolculu%C4%9Funda-sorular-ve-cevaplar-c350463b7ce0> adresinden 13.06.2021'de alınmıştır.
- Şeker, Ş.E. (2013). *İş zekâsı ve veri madenciliği*. İstanbul: Cinius.

- Şimşek, E. ve Canbay, P. (2021). Covid-19 döneminde uzaktan eğitimde mentor gerekliliğinin makine öğrenmesi yaklaşımları ile belirlenmesi ve belirleyicilerin açıklanması. *Avrupa Bilim ve Teknoloji Dergisi*, (26), 246-255.
- Tan, A.L. (2011). Home culture, science, school and science learning: Is reconciliation possible? *Cultural Studies of Science Education*, 6(3), 559–567.
- Terlemez, L. (2008). Eş işlem stratejisi yöntemiyle İMKB'de portföy oluşturmada veri madenciliği uygulaması.
- Thammasiri, D, Delen, D, Meesad, P ve Kasap, N. (2014). A critical assessment of imbalanced class distribution problem: the case of predicting freshmen student attrition. *Expert Systems with Applications*, 41(2), 321–330.
- Toptaş, O. (2021). *Eğitim sektöründe karar destek sistemi ve büyük veri destekli fayda-maliyet analizi*. (Yayımlanmamış Doktora Tezi). İstanbul Kültür Üniversitesi, İstanbul.
- Toptaş, O. ve Şen, A. (2021). Eğitimde ölçme değerlendirme büyük verinin önemi. *Düşünce ve Toplum Sosyal Bilimler Dergisi*, 3(4), 223-243.
- Torgo, L. (2011). *Data mining with R: learning with case studies*. Chapman and Hall/CRC.
- Turgut, M. F. ve Baykul, Y. (2012). *Eğitimde ölçme ve değerlendirme*. Ankara: Pegem Akademi.
- Turing, A. (1950). *Computing machinery and intelligence: Mind*. 59, 433-460.
- Umer, R., Susnjak, T., Mathrani, A. ve Suriadi, S. (2017). On predicting academic performance with process mining in learning analytics. *Journal of Research in Innovative Teaching and Learning*.
- Uzun, P. ve Siemens, G. (2011). Penetrating the Fog: Analytics in Learning and Education. <https://er.educause.edu/articles/2011/9/penetrating-the-fog-analytics-in-learning-and-education> adresinden 16.06.2021'de alınmıştır.
- Ünver, Ö. ve Gamgam, H. (1986). *Uygulamalı istatistik yöntemler*. Ankara: Seçkin Yayıncılık.
- West, Darrell M. (2012). Big data for education: data mining, data analytics, and web dashboards. *Governance Studies at Brookings*, 1-10.
- Witten, D. M., Friedman, J. H. ve Simon, N. (2011). New insights and faster computations for the graphical lasso. *Journal of Computational and Graphical Statistics*, 20(4), 892-900.



- Witten, I. H., Frank, E., Hall, M. A., Pal, C. J. ve Data, M. (2005). *Practical machine learning tools and techniques (3. Baskı)*. In Data Mining.
- Wodecki, A. (2020). *Artificial intelligence in management: self-learning and autonomous systems as key drivers of value creation*. Edward Elgar Publishing.
- Yakut, E. (2012). *Veri madenciliği tekniklerinden C5. 0 algoritması, destek vektör makineleri ile yapay sinir ağlarının sınıflandırma başarılarının karşılaştırılması: İmalat sektöründe bir uygulama*. (Yayımlanmamış Doktora Tezi). Atatürk Üniversitesi Sosyal Bilimleri Enstitüsü, Erzurum.
- Yan, N. ve Au, O. T. S. (2019). Online learning behavior analysis based on machine learning. *Asian Association of Open Universities Journal*.
- Yapay Sinir Ağları Örnek Sorular. (2018). Yapay sinir ağları için örnek sorular. <https://devhunteryz.wordpress.com/2018/05/17/yapay-sinir-aglari-ornek-sorular/> adresinden 12.03.2021’de alınmıştır.
- Yaygın, G. (2019). *Xgboost ve karar ağacı tabanlı algoritmaların diyabet veri setleri üzerine uygulaması*. Yüksek Lisans Tezi, Mimar Sinan Güzel Sanatlar Üniversitesi, Fen Bilimleri Enstitüsü.
- Yekun, E. A., ve Haile, A. T. (2021). Student performance prediction with optimum multilabel ensemble model. *Journal of Intelligent Systems*, 30(1), 511-523.
- Yıldız, M. ve Börekci, C. (2020). Predicting academic achievement with machine learning algorithms. *Journal of Educational Technology and Online Learning*, 3(3), 372-392.
- Yılmaz, H. (2014). *Random forests yönteminde kayıp veri probleminin incelenmesi ve sağlık alanında bir uygulama*. Yüksek Lisans Tezi, Eskişehir Osmangazi Üniversitesi, Eskişehir.
- Yılmaz, M. (2009). Enformasyon ve bilgi kavramları bağlamında enformasyon yönetimi ve bilgi yönetimi. *Ankara Üniversitesi Dil ve Tarih-Coğrafya Fakültesi Dergisi*, 49(1), 95-118.
- Yıldız, M. B., ve Börekci, C. (2020). Predicting academic achievement with machine learning algorithms. *Journal of educational technology and online learning*, 3(3), 372-392.

Zahour, O., Benlahmar, E., Eddaouim, A. ve Hourrane, O. (2020). A comparative study of machine learning methods for automatic classification of academic and vocational guidance questions.

## **EKLER**

### **Ek – 1. Veri Ön İşleme İçin Yazılan Kod**

#Veri ön işleme için kullanılacak fonksiyonlar için gerekli kütüphaneler;

```
Install.packages("funModeling")
```

```
library(funModeling)
```

```
Install.packages("tidyverse")
```

```
library(tidyverse)
```

```
Install.packages("psych")
```

```
library(psych)
```

# Veri seti seçilir.

```
veri <- read.table("data.txt", header=1)
```

# Betimsel istatistiklerin incelenmesi için fonksiyonlar çalıştırılır.

```
summary(veri)
```

```
profiling_num(veri)
```

# Veri setinde kayıp veri olup olmadığı kontrol edilir.

```
sum(is.na(veri))
```

# Korelasyon için ihtiyaç duyulan kod.

```
cor(veri)
```

```
M <- cor(veri)
```

#Görselleştirmek için.

```
corrplot(M, method = "circle")
```

#Yarım görüntü almak için

```
corrplot(M, type = "upper")
```

# Pearson Korelasyon hesaplamak için kullanılır.

```
corr.test(veri, method = "pearson")
```

```
cormatrix <- corr.test(veri, method = "pearson")
```

```
print(cormatrix, short = FALSE)
```

# Sonuları masaüstüne kaydetmek için gerekli olan kod.

```
write.csv(cormatrix$r, "korelasyon.csv")
```

## Ek- 2. C5.0 Karar Ağacı İçin Yazılan Kod

# C5.0 algoritmasının çalıştırılması için C50 kütüphanesi kullanılır.

```
Install.packages("C50")
```

```
library(C5.0)
```

# Veri seti seçilir.

```
veri <- read.table("turkey.dataa.csv", header=1)
```

# Eğitim ve test seti için sample fonksiyonu kullanılarak rastgele seçim yapılır.

```
N = nrow(veri)
```

```
egitim <- sample(1:N, 715, FALSE)
```

# Veri setine ait model tasarlanır.

```
sonuc <- C5.0(instr ~., data = veri[egitim,])
```

# Veri setinin sonuç değerini görselleştirmek için kullanılır.

```
plot(sonuc)
```

# Karar ağaçlarında yaprak düğümleri ile gözlem sayısının ve oranının doğru sınıflandırılıp sınıflandırılmadığı öğrenilir. Bu doğrultuda sınıflandırma kuralı için kullanılan kod;

```
kurallar <- C5.0(instr ~., data = veri[egitim,], rules=TRUE)
```

#Summary fonksiyonu veri seti hakkında özet bilgi sunar.

```
summary(kurallar)
```

# Gerçek değerlerin incelenmesinin ardından karar ağacının eğitim setini ne kadar iyi sınıflandırdığını gözlemlemek amacıyla predict ve table fonksiyonlarının bir kombinasyonu kullanılır.

```
tahmin_egitim <- predict(sonuc, newdata=veri[egitim,], type= "class")
```

```
tahmin_egitim
```

# Eğitim setinin ne kadar iyi sınıflandırma performansı gösterdiği ise table fonksiyonu ile gözlemlenmektedir

```
table(veri$instr[egitim],tahmin_egitim, dnn=c("Observed Class","Predicted Class"))
```

# Test setinin ne kadar iyi sınıflandırma performansı gösterdiği ise table fonksiyonu ile gözlemlenmektedir

```
tahmin_test <- predict(sonuc, newdata = veri[-egitim,], type="class")
```

```
table(veri$instr[-egitim],tahmin_test, dnn=c("Observed Class","Predicted Class"))
```

Eğitim setinin tahmini aşamasında kullanılan kodların test setinde kullanılması için tek bir değişik yapılmaktadır. Bu da eğitim verisinin “-egitim” şeklinde alınmasıdır. R yazılım dilinde “-egitim” demek eğitim seti dışındakiler demektir, bu da test setini ifade etmektedir.

### Ek – 3. Rastgele Orman İçin Yazılan Kod

# Rastgele Orman algoritmasının çalıştırılması için randomForest kütüphanesi çalıştırılır.

```
Install.packages("randomForest")
```

```
library(randomForest)
```

# Veri seti seçilir.

```
veri <- read.table("turkey.dataa.csv", header=1)
```

# Eğitim ve test seti için sample fonksiyonu kullanılarak rastgele seçim yapılır.

```
N <- sample(2, nrow(veri), replace = TRUE, prob = c(0.7, 0.3))
```

```
egitim <- veri[N==1,]
```

```
test <- veri[N==2,]
```

# Bağımsız değişkenler dikkate alınarak mtry değeri belirlenen model kodu.

```
sonuc_egitim <- randomForest(instr ~., data=egitim, ntree = 500, mtry = 32,  
importance = TRUE, proximity = TRUE)
```

```
print(sonuc_egitim)
```

# mtry düzenleme kodu.

```
mtry_düzenleme <- tuneRF(egitim[,-9], egitim[,9],
```

```
stepFactor = 0.5,
```

```
plot = TRUE,
```

```
ntreeTry = 300,
```

```
trace = TRUE,
```

```
improve = 0.05)
```

# mtry düzenlendikten sonra model tekrar çalıştırılır.

```
sonuc_egitim <- randomForest(instr ~., data=egitim, ntree = 500, mtry = 33,  
importance = TRUE, proximity = TRUE)
```

```
print(sonuc_egitim)
```

# Eğitim setinin tahmin performansı için aşağıdaki kod kullanılır.

```
egitim_tahmin <- predict(sonuc_egitim, egitim)
```

```
confusionMatrix(egitim_tahmin, egitim$instr)
```

# Test setinin tahmin performansı için aşağıdaki kod kullanılır.

```
test_tahmin <- predict(sonuc_egitim, test)
```

```
confusionMatrix(test_tahmin, test$instr)
```



#### Ek – 4. Yapay Sinir Ağları İçin Yazılan Kod

# Yapay Sinir Ağları algoritmasının çalıştırılması için neuralnet kütüphanesi çalıştırılır.

```
Install.packages("neuralnet")
```

```
library(neuralnet)
```

# Veri seti seçilir.

```
veri <- read.table("data.txt", header=1)
```

# Normalizasyon ve örneklemin belirlenmesi için gerekli kod aşağıdaki gibidir.  
(Normalizasyon işlemi tüm değişkenlere uygulanmalıdır.)

```
veri$class <- (veri$class - min(veri$class))/(max(veri$class) - min(veri$class))
```

```
veri$nb.repeat <- (veri$nb.repeat - min(veri$nb.repeat))/(max(veri$nb.repeat) -  
min(veri$nb.repeat))
```

```
...
```

```
..
```

```
for(i in 1:33){
```

```
  veri[,i] <- (veri[,i]-min(veri[,i]))/(max(veri[,i]-min(veri[,i])))}
```

# Eğitim ve test seti için aşağıdaki fonksiyon kullanılır.

```
ind <- sample(1:nrow(veri),4074)
```

```
train_data <- veri[ind,]
```

```
test_data <- veri[-ind,]
```

# Veri setine ait model tasarlanır.

```
model <- neuralnet(instr ~., data = training, hidden = 5, err.fct = "ce", linear.output  
= FALSE)
```

```
plot(model)
```

# Çıktı değerinin olasılığı hesaplanır.

```
output <- compute(model , test_data[,-1])
```

```
prediction<-output$net.result*(max(veri[-ind,1])-min(veri[-ind,1]))+min(veri[-ind,1])
```

```
actual <- veri[-ind,1]
```

# Uygun ağırlık değerine ulaşmak için gerçek değer ile tahmin değeri arasındaki minimum fark araştırılır.

```
MSE <- sum((prediction-actual)^2)/nrow(test_data)
```

```
table(actual, round(prediction))
```

# Karışıklık matris sonuçlarına ulaşmak için model değiştirilmeden aşağıdaki kod çalıştırılarak sonuçlara ulaşılır.

```
pred <- neuralnet::compute(model, veri[,c(2:33)])
```

```
pred.2 <- data.frame()
```

```
for(i in 1:5820){
```

```
  pred.2 <- rbind(pred.2, which.max(pred$net.result[i]))}
```

```
pred.2$X1L <- gsub(1, "1", pred.2$X1L)
```

```
pred.2$X1L <- gsub(1, "2", pred.2$X1L)
```

```
pred.2$X1L <- gsub(1, "3", pred.2$X1L)
```

```
prediction <- as.factor(pred.2$X1L)
```

```
reference <- veri[,1]
```

```
confusionMatrix(prediction, reference)
```

## ÖZGEÇMİŞ

### KİŞİSEL BİLGİLER

**Adı ve Soyadı:** Canay CAN

**Doğum Yeri ve Tarihi:**

### EĞİTİM DURUMU

**Lisans Öğrenimi:** Akdeniz Üniversitesi/Eğitim Fakültesi/Matematik ve Fen Bilimleri Eğitimi Bölümü/Fen Bilgisi Öğretmenliği

Atatürk Üniversitesi/Açıköğretim Fakültesi/Çocuk Gelişimi Lisans Programı

**Yüksek Lisans Öğrenimi:** Akdeniz Üniversitesi/Eğitim Bilimleri Enstitüsü/Eğitim Bilimleri Anabilim Dalı/Eğitimde Ölçme ve Değerlendirme Tezli Yüksek Lisans Programı

**Bildiği Yabancı Diller:** İngilizce

**Bilimsel Faaliyetler:**

**Katıldığı Kongreler/Seminerler**

#### 1. Dinleyici

Eğitimde Ölçme Değerlendirme Uygulamaları Ulusal Kongresi, İstanbul, Mayıs, 2021

#### Sertifikalar

20 Ekim 2020 tarihinde Udemy tarafından desteklenen “*R ile Veri Bilimi ve Machine Learning (35 saat)*” adlı çevrimiçi kurs eğitimi.

TÜBİTAK BİDEB 2237-A programı kapsamında desteklenen, Hasan Kalyoncu Üniversitesi Farabi Öğretmen Akademisi bünyesinde 14-19 Haziran 2021 tarihleri arasında düzenlenen “*Lisansüstü Eğitim Gören Araştırmacılara Yönelik Uygulamalı Nicel Veri Analizi Eğitimi*”.

## **İŞ DENEYİMİ**

**Stajlar:** Antalya/Muratpaşa/Barbaros Ortaokulu

**Projeler:** BİDEB – TÜBİTAK, Araştırma Burs ve Destekleri Müdürlüğü, 2209-A Üniversite Öğrencileri Araştırma Projeleri Destekleme Programı (2017)

### **Çalıştığı Kurumlar:**

(2018-2019) Halk Eğitim Merkezi– Geleneksel Türk Okçuluğu Eğiticisi

(2019-2020) Milli Eğitim Bakanlığı– Antalya Kepez Ortaokulu/Fen Bilgisi Öğretmeni

(2021 – devam ediyor) Ekip Network Şirketi– Metin Yazarı

## **İLETİŞİM**

**E-posta Adresi:**

**TARİH:** 11/08/2021

## BİLDİRİM

Hazırladığım tezin/raporun tamamen kendi çalışmam olduğunu ve her alıntıya kaynak gösterdiğimi taahhüt eder, tezimin/raporumun kağıt ve elektronik kopyalarının Akdeniz Üniversitesi Eğitim Bilimleri Enstitüsü arşivlerinde aşağıda belirttiğim koşullarda saklanmasına izin verdiğimi onaylarım:

- Tezimin/Raporumun tamamı her yerden erişime açılabilir.
- Tezim/Raporum sadece Akdeniz Üniversitesi yerleşkelerinden erişime açılabilir.
- Tezimin/Raporumun ..... yıl süreyle erişime açılmasını istemiyorum. Bu sürenin sonunda uzatma için başvuruda bulunmadığım takdirde, tezimin/raporumun tamamı her yerden erişime açılabilir.

11/08/2021

**Canay CAN**

## İNTİHAL RAPORU

tez

ORJİNALLIK RAPORU

% <b>5</b>	% <b>5</b>	% <b>1</b>	%
BENZERLİK ENDEKSİ	İNTERNET KAYNAKLARI	YAYINLAR	ÖĞRENCİ ÖDEVLERİ

BİRİNCİL KAYNAKLAR

<b>1</b>	<a href="http://www.openaccess.hacettepe.edu.tr:8080">www.openaccess.hacettepe.edu.tr:8080</a> İnternet Kaynağı	% <b>1</b>
<b>2</b>	<a href="http://dspace.akdeniz.edu.tr">dspace.akdeniz.edu.tr</a> İnternet Kaynağı	% <b>1</b>
<b>3</b>	<a href="http://egitimbilim.akdeniz.edu.tr">egitimbilim.akdeniz.edu.tr</a> İnternet Kaynağı	<% <b>1</b>
<b>4</b>	<a href="http://dspace.kocaeli.edu.tr:8080">dspace.kocaeli.edu.tr:8080</a> İnternet Kaynağı	<% <b>1</b>
<b>5</b>	<a href="http://dspace.gazi.edu.tr">dspace.gazi.edu.tr</a> İnternet Kaynağı	<% <b>1</b>
<b>6</b>	<a href="http://wiki.socr.umich.edu">wiki.socr.umich.edu</a> İnternet Kaynağı	<% <b>1</b>
<b>7</b>	<a href="http://hakankor.com.tr">hakankor.com.tr</a> İnternet Kaynağı	<% <b>1</b>
<b>8</b>	<a href="http://utek2019.com">utek2019.com</a> İnternet Kaynağı	<% <b>1</b>
<b>9</b>	<a href="http://zoofed.cu.edu.tr">zoofed.cu.edu.tr</a> İnternet Kaynağı	<% <b>1</b>