

**T.C.
HARRAN ÜNİVERSİTESİ
FEN BİLİMLERİ ENSTİTÜSÜ**

YÜKSEK LİSANS TEZİ

**DERİN ÖĞRENME KULLANARAK EL YAZILARINDAN BİLGİ
ÇIKARIMI**

Mehmet TUTAR

BİLGİSAYAR MÜHENDİSLİĞİ ANABİLİM DALI

**ŞANLIURFA
2022**

İÇİNDEKİLER

	Sayfa No
ÖZET	i
ABSTRACT	ii
TEŞEKKÜR.....	iii
ŞEKİLLER DİZİNİ.....	iv
ÇİZELGELER DİZİNİ	v
SİMGELER ve KISALTMALAR DİZİNİ	vi
1. GİRİŞ.....	1
2. ÖNCEKİ ÇALIŞMALAR	5
3. MATERYAL ve YÖNTEM	10
3.1. Materyal.....	12
3.1.1. MFHD Türkçe yazı şablonu hazırlama süreci.....	12
3.1.2. MFHD TTT toplama süreci.....	14
3.1.3. MFHD veri kümesi hazırlık süreci.....	14
3.2. Yöntem	15
3.2.1. Form nitelik bilgisi oluşturma	15
3.2.2. Form hizalama.....	19
3.2.3. Form gürültü temizleme	22
3.2.4. Form karakter algılama, çıkarma ve RGB gri-ton dönüşümü	25
3.2.5. Renkleri ters çevirme	27
3.2.6. Karakter boyut normalleştirme.....	29
3.2.7. Karakter etiketleme	30
4. ARAŞTIRMA BULGULARI ve TARTIŞMA.....	33
4.1. MFHD Üzerinde Yapılan Deneyler.....	34
4.1.1. MFHD-L, MFHD-U ve MFHD-D ile sınıflandırma	35
4.1.2. MFHD'nin niteliklere göre sınıflandırılması	36
4.2. MFHD-D ve MNIST ile Yapılan Deneyler	48
4.3. Tartışma.....	51
5. SONUÇLAR ve ÖNERİLER	53
KAYNAKLAR	55

ÖZET

Yüksek Lisans Tezi

DERİN ÖĞRENME KULLANARAK EL YAZILARINDAN BİLGİ ÇIKARIMI

Mehmet TUTAR

Harran Üniversitesi
Fen Bilimleri Enstitüsü
Bilgisayar Mühendisliği Anabilim Dalı

Danışman: Dr. Öğr. Üyesi Serdar ÇİFTÇİ
YIL: 2022, Sayfa: 60

İnsanlık tarihi boyunca; iletişim kurma, karşılıklı anlaşabilme ve bunun yanında bilginin aktarımı amacıyla çeşitli alfabeler geliştirilmiştir. Teknolojinin ilerlemesi ile birlikte günümüzde el yazısı tanıma için farklı öğrenme yöntemleri geliştirilmiştir. El yazısı tanıma yöntemleri geliştirilirken el yazısı veri kümelerine ihtiyaç duyulmuş ve bütün alfabeler için mümkün olmasa da sık kullanılan alfabeler için el yazısı metin ve karakter veri kümeleri oluşturulmuştur. Bu çalışmada hem el yazısı tanıma hem de el yazısından bilgi çıkarımı için yeni ve çok nitelikli bir veri kümesi sunulmuştur. Bu veri kümesi; Latin harflerden oluşan Türk alfabesi kullanılarak farklı yaş aralığında, farklı eğitim seviyesine sahip, farklı hobileri olan bay ve bayan toplam 20 000 katılımcı tarafından el yazısıyla yazılan, küçük harf (29 Sınıf), büyük harf (29 Sınıf) ve rakam (10 Sınıf) olmak üzere 3 farklı türden ve toplam 68 sınıftan oluşmaktadır. Aynı zamanda Türkçe el yazısı karakter örneklerinin bu ölçekteki kamuya açık ilk veri kümesidir. Küçük harfler 580 000 adet, büyük harfler 580 000 adet ve rakamlar 200 000 adet olmak üzere toplam 1 360 000 adet el yazısı karakter içermektedir. 4 farklı nitelikte (cinsiyet, yaş, eğitim ve hobi) toplanmış; cinsiyete göre 2 (bay, bayan), yaş aralığına göre 4 (5-11 yaş arası, 12-19 yaş arası, 20-30 yaş arası, 31-65 yaş arası), eğitim durumuna göre 4 (ilkokul, ortaokul, lise, yüksekokul) ve hobilere göre 8 (kitap, TV, internet, oyun, spor, müzik, resim, gezi) farklı niteliğe ayrılmış Türkçe el yazısı karakter örneklerinin işlenmiş ve etiketlenmiş ilk veri kümesidir. Sunduğumuz veri kümesi, sadece bilgisayar bilimleri alanında değil farklı bilimsel alanlarda çalışma yapan araştırmacılara da katkı verecek niteliktedir. Türkiye’de en geniş katılımcı sayısı ile toplanan el yazısı veri kümesinde bulunan küçük harf, büyük harf ve rakamlara, el yazısı karakter tanıma için bilinen sınıflandırma algoritmaları uygulanmış ve bu yöntemlerin performansları incelenmiştir. Bu çalışmada el yazısı bilgisi ile cinsiyet, eğitim düzeyi, hobi ve yaş grubu bilgisi arasında bir ilişki olup olmadığı araştırılmıştır. Veri kümesi ile yapılan deney sonuçları incelendiğinde; performans başarımları sırasıyla cinsiyet, eğitim durumu, yaş grubu ve hobi durumuna göre olmuştur. Buradan hareketle kişinin cinsiyet bilgisinin el yazısına daha fazla yansıdığı gözlemlenmiştir. Sunmuş olduğumuz veri kümesinde bulunan el yazısı rakamlarla yapılan sınıflandırma performansı, kıyaslama veri kümesi olan MNIST ile karşılaştırılmış ve sonuçları tartışılmıştır.

ANAHTAR KELİMELER: MFHD, El Yazısı Karakter Tanıma, Çok özellikli OCR Veri Kümesi, El Yazısı İnsan-Karakter İlişkisi, El Yazısı Fizyolojik İlişkiler

ABSTRACT

MSc Thesis

EXTRACTING INFORMATION FROM HANDWRITING USING DEEP LEARNING

Mehmet TUTAR

Harran University
Graduate School of Natural and Applied Sciences
Department of Computer Engineering

Supervisor: Assist. Prof. Dr. Serdar İFTÇİ
Year: 2022, Page: 60

Throughout human history, various alphabets have been developed for communication, mutual understanding, and knowledge transferring. With the advancement of technology, different learning methods have been developed for handwriting recognition. While developing handwriting recognition methods, handwritten datasets were needed, and although it is not possible for all alphabets, some frequently used alphabet's handwritten text and character datasets were created. This study presents a new and highly qualified dataset for handwriting recognition and information extraction from handwriting which consists of three different types with lowercase letters (29 classes), uppercase letters (29 classes), and digits (29 classes) with a total of 68 classes, written by a total of 20 000 unique participants of varying genders, hobbies, ages, and education levels by using the Turkish alphabet consisting of Latin letters. The dataset is also the first publicly available dataset of Turkish handwritten character samples with this scale. It contains 1 360 000 handwritten characters, 580 000 lowercase letters, 580 000 uppercase letters, and 200 000 digits. It is the first processed and labeled Turkish handwritten dataset that collected with four different features (gender, age, education, and hobby), and it varies for genders (male, female), age groups (5 to 11, 12 to 19, 20 to 30, and 31 to 65 age), an education level (elementary school, middle school, high school, college+), and hobbies (books, TV, internet, games, sports, music, painting, travel). The dataset we offer not only contributes to researchers working in computer science but also in different scientific fields. It is the largest handwritten dataset collected from a considerable number of participants, and known handwritten classification algorithms were applied to the lowercase/uppercase letters and digits, and their performances were examined. This study investigates the relationship between handwriting and knowledge of gender, education level, hobby, and age group. When the results of the experiments conducted on our dataset were examined, the performance rate ranked, namely gender, education level, age group, and hobby status. From this point of view, it has been observed that the gender of a person is more reflected in handwriting. The classification performance of handwritten digits in the dataset we have presented is compared with the benchmark dataset MNIST, and their results are discussed.

KEYWORDS: MFHD, Handwritten Character Recognition, Multi-featured OCR Dataset, Handwritten Human-Character Relationship, Handwritten Physiological Relationships

TEŞEKKÜR

Bu araştırmanın konusu, veri kümesinin hazırlanması, deneysel çalışmaların yapılması, sonuçların değerlendirilmesi ve tezin yazımı aşamasında, gece gündüz desteğini esirgemeyen ve Lisansüstü eğitimim süresince birçok konuda bilgi ve birikimleri ile destek olan, birlikte çalışmaktan onur duyduğum değerli hocam ve tez danışmanım Sayın Dr. Öğr. Üyesi Serdar ÇİFTÇİ' ye çok teşekkür ederim.

TTT (Türkçe Yazı Şablonu) ile el yazısı karakter formlarının toplanması sürecine katkılarından dolayı Millî Eğitim Bakanlığı (MEB) ve Harran Üniversitesine teşekkür ederim. Ayrıca veri kümesinin toplanması ve hazırlanması aşamasında destek olan, yardımlarını esirgemeyen formların doldurulduğu bütün şehir ve üniversitelerde bulunan hocalarımıza, MEB'de görev yapan öğretmenlerimize ve tüm katılımcılara teşekkür ederim.

Doğduğum günden itibaren şahsıma, her koşulda destek veren ve maddi manevi desteklerini hiçbir zaman esirgemeyen değerli annem ve babama şükranlarımı sunarım. Bu süreçte her zaman destek olan ve fedakârlık yapan, her an göstermiş olduğu sabır ve anlayışı için sevgili eşime çok teşekkür ederim. Ayrıca bana her zaman umut olan, süreç içinde güzel hatıralar yaşadığımız ve ihmal ettiğim sevgili oğlum ve kızlarıma sabır ve anlayışları için çok teşekkür ederim.

ŞEKİLLER DİZİNİ

Sayfa No

Şekil 3.1 Formlarda bulunan nitelik sayılarının dağılımı, (a) cinsiyete göre (b) eğitim durumuna göre (c) yaş grubuna göre (d) hobi durumuna göre	11
Şekil 3.2 TTT formlarından rastgele seçilen küçük harf pangram alanı	12
Şekil 3.3 Doldurulmuş TTT formlarından rastgele seçilen bir örnek (hizalanmamış)	13
Şekil 3.4 MFHD-Feature dosya içeriği	15
Şekil 3.5 Katılımcıların nitelik yaş dağılımları. (a) cinsiyete göre (b) eğitim durumuna göre (c) hobi durumuna göre.....	17
Şekil 3.6 MFHD veri kümesinin cinsiyet, eğitim durumu ve yaş grubuna göre dağılım grafiği.....	18
Şekil 3.7 Eksen kayması bulunan rastgele örnek bir form (En/boy oranı korunmuş)	20
Şekil 3.8 Eksen kayması olan hizalanmış aynı form (Bakınız Şekil 3.7).....	21
Şekil 3.9 Rastgele seçilen gürültülü bir form (hizalanmamış)	23
Şekil 3.10 Gürültüsü giderilmiş aynı (Bakınız Şekil 3.9) form (hizalanmış)	24
Şekil 3.11 Formlardan çıkarılan rastgele küçük harf (29 sınıf) karakter örnekleri.....	26
Şekil 3.12 Formlardan çıkarılan rastgele büyük harf (29 sınıf) karakter örnekleri	26
Şekil 3.13 Formlardan çıkarılan rastgele rakam (10 sınıf) karakter örnekleri.....	27
Şekil 3.14 Renkleri ters çevrilen rastgele küçük harf (29 sınıf) karakter örnekleri	28
Şekil 3.15 Renkleri ters çevrilen rastgele büyük harf (29 sınıf) karakter örnekleri.....	28
Şekil 3.16 Renkleri ters çevrilen rastgele rakam (10 sınıf) karakter örnekleri	29
Şekil 3.17 Boyutu normalleştirilen rastgele karakter örnekleri (28x28 piksel boyutunda)	29
Şekil 3.18 Etiketlenmiş rastgele küçük harf karakter örnekleri.....	30
Şekil 3.19 Etiketlenmiş rastgele büyük harf karakter örnekleri.	31
Şekil 3.20 Etiketlenmiş rastgele rakam karakter örnekleri.....	32
Şekil 4.1 KNN ve LeNet5 yöntemleri ile MFHD (MFHD-L, MFHD-U, MFHD-D) sınıflandırma sonuçları	36
Şekil 4.2 Cinsiyet ve eğitim durumuna göre elde edilen ortalama sınıflandırma doğruluk sonuçlarının (Çizelge 4.2, Çizelge 4.5 ve Çizelge 4.8) sıcaklık haritası gösterimi. Kullanılan kısaltmalar: MA: Bay, FE: Bayan, ES: İlkokul, MS: Ortaokul, HS: Lise, C+: Yüksekokul, K5: KNN (k=5), L5: LeNet5.....	46
Şekil 4.3 Yaş grubuna göre elde edilen ortalama sınıflandırma doğruluk sonuçlarının (Çizelge 4.3, Çizelge 4.6 ve Çizelge 4.9) sıcaklık haritası gösterimi. Kullanılan kısaltmalar: G1: 5-11 yaş arası, G2: 12-19 yaş arası, G3: 20-30 yaş arası, G4: 31-65 yaş arası, K5: KNN (k=5), L5: LeNet5	46
Şekil 4.4 Hobi durumuna göre elde edilen ortalama sınıflandırma doğruluk sonuçlarının (Çizelge 4.4, Çizelge 4.7 ve Çizelge 4.10) sıcaklık haritası gösterimi. Kullanılan kısaltmalar: BO: Kitap, TV: Televizyon, IN: İnternet, GA: Oyun, SP: Spor, MU: Müzik, PA: Resim, TR: Gezi, K5: KNN (k=5), L5: LeNet5	47
Şekil 4.5 KNN yöntemi ile MNIST ve MFHD-D sınıflandırma sonuçları	49
Şekil 4.6 Farklı modeller ile MNIST ve MFHD-D sınıflandırma sonuçları	50

ÇİZELGELER DİZİNİ

Sayfa No

Çizelge 2.1 Literatürde bulunan bazı veri kümelerine ait bilgiler. Kullanılan kısaltmalar: (c): karakter, (w): kelime anlamını ifade etmektedir	9
Çizelge 4.1 KNN ve LeNet5 yöntemleri ile MFHD (MFHD-L, MFHD-U, MFHD-D) sınıflandırma sonuçları	35
Çizelge 4.2 MFHD-L veri kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak cinsiyet ve eğitim durumuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama	38
Çizelge 4.3 MFHD-LA örneklem kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak yaş grubuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama	39
Çizelge 4.4 MFHD-LH örneklem kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak hobi durumuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama	40
Çizelge 4.5 MFHD-U veri kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak cinsiyet ve eğitim durumuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama	41
Çizelge 4.6 MFHD-UA örneklem kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak yaş grubuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama	42
Çizelge 4.7 MFHD-UH örneklem kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak hobi durumuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama	43
Çizelge 4.8 MFHD-D veri kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak cinsiyet ve eğitim durumuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama	44
Çizelge 4.9 MFHD-DA örneklem kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak yaş grubuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama	44
Çizelge 4.10 MFHD-DH örneklem kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak hobi durumuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama	45
Çizelge 4.11 KNN yöntemi ile MNIST ve MFHD-D sınıflandırma sonuçları	49
Çizelge 4.12 Farklı modeller ile MNIST ve MFHD-D sınıflandırma sonuçları	50

SİMGELER ve KISALTMALAR DİZİNİ

ANN	Yapay Sinir Ağları
AVG	Ortalama
BO	Kitap
CNN	Evrışimli Sinir Ağları
CSV	Virgülle Ayrılmış Değerler
C+	Yüksekokul
DL	Derin Öğrenme
DPI	İnç Başına Düşen Nokta Sayısı
EMNIST	Genişletilmiş MNIST
ES	İlkokul
FE	Bayan
G1	5-11 Yaş Arası
G2	12-19 Yaş Arası
G3	20-30 Yaş Arası
G4	31-65 Yaş Arası
GA	Oyun
HCR	El yazısı Karakter Tanıma
HS	Lise
IN	İnternet
K5	KNN k=5
KNN	K-En yakın Komşu
KR	Karakter
KVKK	Kişisel Verileri Koruma Kanunu
L5	LeNet5
MA	Bay
MFHD	Çok Özellikli El Yazısı Veri Kümesi
MFHD-D	MFHD Rakam
MFHD-DA	MFHD-D Yaş
MFHD-DH	MFHD-D Hobi
MFHD-L	MFHD Küçük Harf
MFHD-LA	MFHD-L Yaş
MFHD-LH	MFHD-L Hobi
MFHD-U	MFHD Büyük Harf
MFHD-UA	MFHD-U Yaş
MFHD-UH	MFHD-U Hobi
ML	Makine Öğrenme
MLP	Çok Katmanlı Algılayıcı
MNIST	Biçimlendirilmiş NIST
MS	Ortaokul
MU	Müzik
NIST	ABD Ulusal Standartlar ve Teknoloji Enstitüsü
OCR	Optik Karakter Tanıma
PA	Resim
RGB	Kırmızı Yeşil Mavi
SIFT	Ölçek Değişmez Özellik Dönüşümü
SP	Spor
SURF	Hızlandırılmış Sağlam Özellik
SVM	Destek Vektör Makinesi
TR	Gezi
TTT	Türkçe Yazı Şablonu
TV	Televizyon

1. GİRİŞ

Karakter tanıma araştırması alanındaki ilk önemli girişim 1959 yılında Grimsdale ve ark. tarafından yapılmıştır (Grimsdale ve ark., 1959). 1995'te NIST (ABD Ulusal Standartlar ve Teknoloji Enstitüsü), el yazısı formlardan bölümlere ayrılmış karakterlerin (harfler ve rakamlar) ikili görüntülerini içeren ve 62 etiketli sınıftan oluşan "Özel Veri Kümesi 19" veri kümesini yayımlamıştır (Grother, 1995). 1998'de yalnızca rakamlardan (60 000 eğitim, 10 000 test örneği) oluşan, 28x28 piksel boyutunda, ön işleme tabi tutulmuş ve biçimlendirilmiş MNIST (Modified-NIST) veri kümesi oluşturulmuş ve bu veri kümesi, rakam tanıma çalışmaları için bir referans haline gelmiştir (Lecun ve ark., 1998). NIST veri setinin 2. versiyonu 2017 yılında 3 699 el yazısı örnek formunun tam sayfa ikili görüntüleri ve aynı 62 sınıfın 814 255 örnek rakam ve karakteri ile yayımlanmıştır (Grother, 2017). Karakter tanıma için NIST veri kümesi bir referans haline gelmiştir. 2017 yılında yayımlanan genişletilmiş MNIST (Extended-MNIST), MNIST veri kümesi ile aynı formatta düzenlenmiş, büyük harf, küçük harf ve rakam görüntülerinden oluşan daha kapsamlı bir veri kümesidir (Cohen ve ark., 2017).

Daha önce yayımlanan veri kümelerinden yapısal olarak farklı olan el yazısı karakter veri kümesi C-Cube (Cursive Character Challenge) yayımlanmıştır (Camastra ve ark., 2006). C-Cube veri seti; el yazısı karakterlerin bit haritası (bitmap) ile saklandığı, Latin harflerin 26 küçük ve 26 büyük versiyonunu içeren 57 293 karakterden oluşmaktadır.

Teknoloji gelişimine bağlı olarak el yazısı analizinde de önemli aşamalar kaydedilmiştir (Demir ve Uğurlu, 2021). El yazısından kişiye ait bilgi çıkarımı ve farklı tanımlayıcı bilgiler elde etme isteği insanlara ilgi çekici gelmiştir. Ünlü tarihçilerden birisi olarak kabul edilen Suetonius Tranquillus, imparatorların farklı el yazısı yazım şekline sahip olduklarını fark eden ve bu konuyu ilk analiz eden kişi olarak bilinmektedir (Roberts, 2002). Tranquillus, farklı el yazısı yazım şekillerinin farklı karakteristik özellikleri yansıttığına dair düşüncesi ilk ifade edendir. Aynı

şekilde 17. yy. da yaşayan yazarlar da farklı el yazısı yazım biçimleri ile kişinin karakteristik özelliklerini ilişkilendirmeye çalışmışlardır (Uğurlu ve ark., 2010). El yazısı, kişilerin durumlarından ve farklı özelliklerinden kolayca etkilenebilmektedir (Birincioğlu ve ark., 2010). DNA örneği, kan grubu ve parmak izi gibi bireysel farklılık gösteren başlıca kişisel özelliklerdendir. Bu sebeple de el yazısı, geçerliliği kabul görmüş ayırt edici kişisel bir özelliktir (Huber ve Headrick, 1999).

Var olan veri kümeleri, el yazısı karakter tanıma modellerinin eğitimleri için kullanılmaktadır. Ancak tüm alfabeler ve diller için bu veri kümeleri yeterli olmamaktadır. Çünkü her alfabe ve dilin kendine özgü karakter yapısı ve şekli bulunmaktadır. Bu kapsamda araştırmacıların farklı alfabe ve dillerde el yazısı tanıma sistemi geliştirmesi için yeni veri kümelerine ihtiyaç duyulmaktadır. Bazı diller için mevcut veri kümeleri bulunmasına karşın Türkçe için kapsamlı bir veri kümesi bulunmaması, araştırmacılar için zorluk oluşturmakta ve yeterince araştırma konusu haline gelememektedir.

Türkçe el yazısı karakterlerden oluşan kapsamlı ve nitelikli bir veri kümesi oluşturmamızın diğer motivasyonu, mevcutta var olan bazı veri kümelerinin yalnızca temel Latin alfabesinin karakterlerini içermesi ve ç, ğ, ı, ö, ş, ü gibi noktalı ve çentikli harfleri içermemesidir. Farklı harflerin bulunduğu alfabelerde iyi performans gösteren el yazısı karakter tanıma yöntemleri geliştirebilmek için o harflerin bulunduğu el yazısı karakter veri kümelerine ihtiyaç vardır. Bu geçerli sebepler göz önüne alındığında Türkçe el yazısı karakter örneklerinin bulunduğu veri kümesinin hazırlanması önem arz etmektedir.

Yaptığımız araştırmalar neticesinde, karakter sayısı bakımından bu ölçekte bir Türkçe el yazısı karakter veri kümesi mevcut değildir. Bu çalışmamızla, bahsi geçen bu eksikliği gidermek için kapsamlı bir Türkçe el yazısı karakter veri kümesi hazırlanmıştır.

Bu çalışma ile sunmuş olduğumuz MFHD (Çok Özellikli El Yazısı Veri Kümesi); El yazısı karakter tanıma (HCR) uygulamaları için kullanılan, katılımcı

sayısı bakımından en kapsamlı ve nitelik bilgisi içeren ilk veri kümesidir. Aynı zamanda el yazısı karakter sayısı bakımından bu büyüklükte Türkçe el yazısı karakter örneklerini içeren kamuya açık ilk veri kümesidir.

Bu makalede, Latin harflerden oluşan Türk alfabesi kullanılarak farklı yaş aralığında, farklı eğitim seviyesine sahip, farklı hobileri olan bay ve bayan toplam 20 000 katılımcı tarafından el yazısıyla yazılan; küçük harf (29 Sınıf), büyük harf (29 Sınıf) ve rakam (10 Sınıf) olmak üzere 3 farklı türden ve toplam 68 sınıftan oluşan kamuya açık bir veri kümesi sunulmuştur. Bu veri kümesi 1 360 000 el yazısı karakter görüntüsünden oluşmaktadır.

Bu çalışmayla literatüre sunmayı amaçladığımız katkılar aşağıda sıralanmıştır:

- 20 000 farklı katılımcı tarafından yazılan, 4 farklı nitelik (Feature) bilgisi içeren formlardan; toplanan, işlenen ve etiketlenen dünyada bilinen en büyük çevrimdışı (offline) el yazısı karakter veri kümesidir.
- Türkiye’de en geniş katılımcı sayısı ile toplanan ve görüntü sayısı bakımından Türkçe el yazısı karakter örneklerinin kamuya açık ilk veri kümesidir.
- Türkçe el yazısı karakter örnekleri bakımından büyük ölçekli ilk veri kümesidir. Küçük harfler 580 000 adet, büyük harfler 580 000 adet ve rakamlar 200 000 adet olmak üzere toplam 1 360 000 adet el yazısı karakter içermektedir.
- 4 farklı nitelikte (cinsiyet, yaş, eğitim ve hobi) toplanmış; cinsiyete göre 2 (bay, bayan), yaş aralığına göre 4 (5-11 yaş, 12-19 yaş, 20-30 yaş, 31-65 yaş), eğitim durumuna göre 4 (ilkokul, ortaokul, lise, yüksekokul) ve hobilere göre 8 (kitap, TV, internet, oyun, spor, müzik, resim, gezi) farklı niteliğe ayrılmış Türkçe el yazısı karakter örneklerinin işlenmiş ve etiketlenmiş ilk veri kümesidir.

- Latin harfler kullanılarak el yazısı ile yazılan rakam karakter sayısı bakımından toplanan, işlenen ve etiketlenen en büyük veri kümesidir.
- Katılımcı sayısı bakımından, el yazısıyla insan karakterleri ve fizyolojik yapıları arasında bir ilişki olup olmadığının yapay zekâ ve derin öğrenme yöntemleri ile araştırıldığı ve sonuçlarının gözlemlendiği ilk çalışmadır.
- Sunduğumuz veri kümesi, sadece bilgisayar ve mühendislik bilimleri alanında değil farklı bilimsel alanlarda (Grafoloji, Sosyal ve Beşerî bilimler) çalışma yapan araştırmacılara da katkı sunacak niteliktedir.
- Bu veri kümesi kullanılarak mevcutta (state-of-art) var olan el yazısı karakter sınıflandırma yöntemlerinin performansları karşılaştırılmış ve o yöntemlerin tutarlılığı incelenmiştir.
- MFHD veri kümesi, bilimsel çalışmalara katkıda bulunmak üzere <https://github.com/TezKabulAldığında/MFHD.git> adresinde araştırmacıların erişimine sunulacaktır.

Çalışmanın geri kalanı aşağıdaki gibi yapılandırılmıştır. Bölüm II'de Önceki çalışmalardan bahsedilmiştir. Bölüm III'te materyal ve yöntemler açıkça belirtilmiş ve MFHD veri kümesinin oluşturulması ayrıntılı olarak açıklanmıştır. Bölüm IV'te araştırma ve bulgular detaylı olarak ele alınarak MFHD veri kümesinin state-of-art yöntemler üzerindeki performansı incelenmiş ve MNIST veri kümesi ile tutarlılığı analiz edilmiştir. Bölüm V'de sonuçlar ve gelecekteki öneriler üzerinde durularak MFHD veri kümesinin el yazısından bilgi çıkarma ve el yazısı karakter tanıma konusunda literatüre olan katkısına yer verilmiştir.

2. ÖNCEKİ ÇALIŞMALAR

El yazısı sembolleri ve karakterleri bilgisayar sistemi tarafından tanıma işlemine el yazısı tanıma denir (Bhatia, 2014). Geçmiş çalışmalarda, makine okumasının insan okumasından performans olarak önemli ölçüde düşük olduğu gösterilmiştir (Arica ve Yarman-Vural, 2001). Ancak günümüzde gelişen donanım mimarisi ve artan veri kümelerine bağlı olarak bu durum değişmiştir.

El yazısı tanıma, çevrimdışı (etkileşimsiz) ve çevrimiçi (etkileşimli) yöntemler olarak sınıflandırılabilir (Plamondon ve Srihari, 2000). Literatürde çevrimdışı (daha önceden kâğıt üzerine yazılmış bilgilerin dijitalleştirilerek tanıma çalışması) ve çevrimiçi (el yazısının yazıldığı esnada tanımaya çalışılması) el yazısı tanıma olarak bilinen birçok araştırma bulunmaktadır.

El yazısı karakter tanıma çalışmaları için literatürdeki bazı farklı dil ve alfabelerden oluşturulan Arapça (Al-Ohali ve ark., 2003; El-Sherif ve Abdelazeem, 2007; Mahmoud ve ark., 2014), Çince (Saito, 1985; Liu ve ark., 2011), Hintçe (Manjusha ve ark., 2019), Korece (Kim ve ark., 1996), Fransızca (Arvanitopoulos ve ark., 2017), İspanyolca (Toselli ve ark., 2004; Espana ve ark., 2004), Rusça (Nurseitov ve ark., 2021) ve Latince (Grother, 1995; Marti ve Bunke, 2002; Bartos ve ark., 2020) gibi çeşitli karakter veri kümeleri (Benchmark datasets) bulunmaktadır.

Literatürde el yazısı karakter veri kümelerinin yanında farklı dil ve alfabeler kullanılarak sadece rakamlardan oluşan Latince (Lecun, 1998; Srl, 1994; Hull, 1994; De Campos ve ark., 2009; Kusetogullari ve ark., 2020), Çince (Liu ve ark., 2013), Farsça (Mohammad ve ark., 2011) ve Arapça (El-Sherif ve Abdelazeem, 2007) gibi çeşitli rakam veri kümeleri de bulunmaktadır.

Literatür incelememizde Türkçe el yazısı karakter tanıma ve geliştirme ile ilgili çok az çalışmaya rastladık. Türkçe el yazısı tanıma için yapılan bazı çalışmalarda

araştırmacılar kamuya açık olmayan kendilerine ait veri kümelerini kullanmışlardır (Vural ve ark., 2004; Şekerci ve Kandemir, 2006; Şekerci, 2007; Bartos ve ark., 2018).

2003 yılında yaklaşık 20 000 karakter ile Türkçe büyük harflerden oluşan ilk el yazısı karakter veri kümesi oluşturulmuştur (Çapar ve ark., 2003). Türkçe el yazısı karakter tanıma üzerine yapılan en kapsamlı veri kümesi çalışması 2020 yılında T-H-E veri kümesi adıyla yapılmıştır (Bartos ve ark., 2020). Bu veri kümesi; 200 katılımcıdan toplanan el yazısı Türkçe, Macarca ve İngilizce karakter içeren, 78 farklı sınıftan oluşan, küçük harf ve büyük harf dâhil 156 000 karakter içeren kamuya açık bir veri kümesidir.

Literatür araştırmamızda; farklı alfabe ve diller kullanılarak oluşturulan el yazısı veri kümeleri üzerinde, birçok metot ve yöntemle yapılan el yazısı karakter tanıma ve sınıflandırma çalışmalarına rastladık. Literatürde kullanılan bazı yöntem ve sınıflandırma çalışmalarına tezimizde değinilmiştir.

K-en yakın komşu (KNN: K-Nearest Neighbor) algoritması ilk olarak 1952'de araştırmacılar tarafından yayımlandı (Fix ve Hodges, 1952) ve 1967'de yeniden düzenlendi (Cover ve Hart, 1967). KNN, performansı kullanılan mesafe metriğine bağlı etkili bir veri sınıflandırma yöntemidir (Wang ve ark., 2018; Maillou ve ark., 2015).

KNN, sınıflandırma problemlerini etkin bir şekilde çözmek için kullanılabilir başarılı bir denetimli makine öğrenme algoritmasıdır (Arslan ve Arslan, 2021). Uygun metrik değerlerini belirlemek için KNN üzerine çeşitli araştırma çalışmaları yapılmıştır (Cha, 2007; Abu Alfeilat ve ark., 2019).

Zhang ve ark. (2017), sabit bir k değeri kullanmak yerine KNN sınıflandırmasına bir eğitim aşaması daha ekleyerek, farklı optimal k değerlerini öğrenen bir kTree yöntemi ile daha yüksek sınıflandırma doğruluğu elde etmişlerdir.

Grover ve Toghi (2019), MNIST veri kümesi üzerinde K-En Yakın Komşu (KNN) mesafe metriğini kayan pencere tekniği ile değiştirerek sınıflandırma doğruluk oranını iyileştirmişlerdir.

MNIST veri kümesi üzerinde Evrişimli Sinir Ağları (CNN) kullanılarak yapılan çalışmalarda başarılı sonuçlar elde edildiği gözlemlenmiştir (Wu, 2018; Bharadwaj ve ark., 2020; El-Sawy ve ark., 2017). CNN ve Yapay Sinir Ağları (ANN) kullanılarak MNIST veri seti üzerinde yapılan çalışmada, CNN ağının, görüntü sınıflandırma için ANN ağına göre daha iyi performans ortaya koyduğu belirtilmektedir (Beohar ve Rasool, 2021).

Belirgin yerel nitelikleri daha iyi çıkarmak ve el yazısı görüntülerdeki gürültüyü filtrelemek için dikkat mekanizması tabanlı bir sinir ağı önerilmiştir (Hao ve Chen, 2020). MNIST veri kümesi üzerinde yapılan bu çalışmada, LeNet5 (Lecun ve ark., 1998) ağına dikkat mekanizması eklenerek oluşturulan modelin doğruluk oranını iyileştirdiği bildirilmektedir.

Destek Vektör Makinesi (SVM: Support Vector Machine), Vapnik tarafından ortaya çıkarılan ve makine öğrenme uygulamalarında fazla ilgi gören bir yöntemdir (Vapnik, 1995). Destek vektör makinesinin temel dayanağı, istatistiksel öğrenme teorisi başka bir deyişle Vapnik-Chervonenkis (VC) teorisidir (Li ve ark., 2009).

SVM; yüz tanıma ve doğrulama, konuşmacı tanıma, metin sınıflandırma, tahmin yürütme, el yazısı karakter ve metin tanıma gibi birçok farklı alanda başarıyla uygulanmıştır (Oliveira ve Sabourin, 2004). Son zamanlarda yapılan birçok araştırma, SVM'nin genellikle, diğer veri sınıflandırma algoritmalarından daha iyi ayırt edicilik göstererek iyi performans ortaya koyabildiğini göstermiştir (Hearst ve ark., 1998).

MNIST veri kümesi üzerinde tüm makine öğrenme algoritmalarının uygulandığı çalışmada, el yazısıyla yazılmış rakamların sınıflandırılmasında SVM

yöntemi ile en yüksek doğruluk oranına ulaşıldığı ifade edilmektedir (Gope ve ark., 2021).

MNIST veri kümesi üzerinde, Makine öğrenme (ML) yöntemlerinin, Derin öğrenme (DL) modelleri ile birlikte kullanılmasıyla (CNN-SVM hibrit modeli) doğruluk oranının iyileştirildiği bildirilmektedir (Ahlawat ve Choudhary, 2020; Yu ve ark., 2015).

Demirkaya ve Çavuşoğlu (2021); el yazısı karakter tanımadaki sık kullanılan ML ve DL algoritmalarının performansını karşılaştırarak en başarılı modelin CNN olduğunu göstermişlerdir. MNIST veri kümesi üzerinde, ML ve DL yöntemlerinden olan SVM, Çok Katmanlı Algılayıcı (Multi-layer Perceptron (MLP)) ve CNN modelleri uygulanmış ve CNN'nin daha iyi sonuçlar verdiği ifade edilmiştir (Pashine ve ark., 2021).

El yazısı karakter tanımadaki performans başarımı, büyük verilerle çalışabilmesi ve karmaşık mimarilere çözüm üreten donanım gelişimine bağlıdır (Baldominos ve ark., 2019). An ve ark. (2020), yüksek performansın tek bir yöntem veya algoritma ile elde edilemeyeceğini, birçok yöntemin bir arada kullanılmasıyla elde edilebileceğini göstermişlerdir. Bu çalışma, eğitim setinden bağımsız olarak eğitilmiş üç modeli kullanan ve bu modellerden gelen sonuçları çoğunluk oylamasına tabi tutarak karar veren bir yöntem içermektedir.

Performans artırımı için sınıflandırma seçiminin yanı sıra özellik çıkarma tekniklerinin de uygun belirlenmesi gerekmektedir (Purohit ve Chauhan, 2016). Bununla birlikte veri kümesinin büyüklüğü ve karakterlerin benzersizliği de performansı artırmak için önemli bir kriterdir. Literatürde bulunan bazı veri kümelerine ait bilgiler Çizelge 2.1'de gösterilmiştir.

Çizelge 2.1 Literatürde bulunan bazı veri kümelerine ait bilgiler. Kullanılan kısaltmalar: (Ka): karakter, (Ke): kelime anlamını ifade etmektedir

Veri Kümeleri	Küçük Harf	Büyük Harf	Rakam	Katılımcı Sayısı	Toplam Karakter / Kelime	Cinsiyet	Eğitim	Yaş	Ho bi
1 Çapar ve ark. (2003)	X	✓	✓	X	27 000 (Ka)	X	X	X	X
2 Şekerci ve Kandemir (2006)	✓	✓	X	172	9 976 (Ka)	X	X	X	X
3 NIST SD 19 (Grother, 2017)	✓	✓	✓	3 600	814 255 (Ka)	X	X	X	X
4 MNIST (Lecun, 1998)	X	X	✓	500	70 000 (Ka)	X	X	X	X
5 EMNIST (Cohen ve ark., 2017)	✓	✓	✓	3 600	814 255 (Ka)	X	X	X	X
6 C-Cube Database (Camastra ve ark., 2006)	✓	✓	X	X	57 293 (Ka)	X	X	X	X
7 Hasyv2 Dataset (Thoma, 2017)	✓	✓	✓	477	168 233 (Ka)	X	X	X	X
8 T-H-E Dataset (Bartos ve ark., 2020)	✓	✓	X	200	156 000 (Ka)	X	X	X	X
9 Şekerci (2007)	✓	✓	✓	172	11 696 (Ka)	X	X	X	X
10 CASIA Database (Liu ve ark., 2011)		Belirtilmemiş		1 020	1 350 000 (Ka)	X	X	X	X
11 PE92 Korean Database (Kim ve ark., 1996)		Belirtilmemiş		500	235 000 (Ka)	X	X	X	X
12 AHDBase Database (El-Sherif ve Abdelazeem, 2007)	X	X	✓	700	70 000 (Ka)	X	X	X	X
13 KHATT Database ¹ (Mahmoud ve ark., 2014)		Kelime		1 000	178 255 (Ke)	✓	✓	✓	X
14 IAM Database (Marti ve Bunke, 2002)		Kelime		400	82 227 (Ke)	X	X	X	X
15 Malamayam Database ² (Manjusha ve ark., 2019)		Belirtilmemiş		77	29 302 (Ka)	✓	X	✓	X
16 Toselli ve ark. (Toselli ve ark., 2004)		Kelime		29	2 127 (Ke)	X	X	X	X
17 The SPARTACUS-Database (España ve ark., 2004)		Kelime		1 500	100 000 (Ke)	X	X	X	X
18 CFRAMUZ Dataset (Arvanitopoulos ve ark., 2017)		Kelime		1	18 000 (Ke)	X	X	X	X
19 HKR Database (Nurseitov ve ark., 2021)		Belirtilmemiş		200	715 699 (Ka)	X	X	X	X
20 MFHD (Bizimki)	✓	✓	✓	20 000	1 360 000 (Ka)	✓	✓	✓	✓

¹ Bu çalışmada katılımcılardan cinsiyet, eğitim ve yaş grubu bilgisi alınmış fakat veri kümesi ile birlikte paylaşılmamıştır

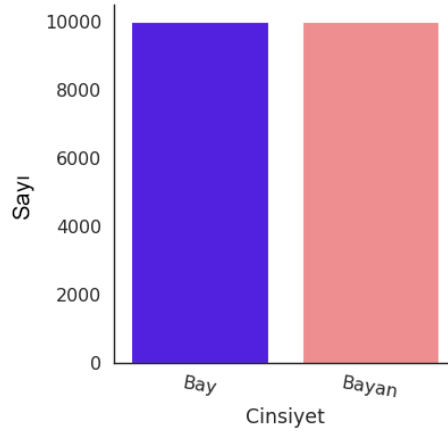
² Bu çalışmada katılımcılardan cinsiyet ve yaş grubu bilgisi alınmış fakat veri kümesi ile birlikte paylaşılmamıştır

3. MATERYAL ve YÖNTEM

MFHD (Çok Özellikli El Yazısı Veri Kümesi / Multi-Featured Handwritten Dataset), 4 farklı nitelik türü (cinsiyet, yaş, eğitim, hobi) içeren, 20 000 tekil katılımcıdan toplanmış kamuya açık Türkçe el yazısı veri kümesidir. Bu veri kümesi Türk alfabesinde bulunan küçük harf (29 sınıf), büyük harf (29 sınıf) ve rakamlar (10 sınıf) olmak üzere 3 farklı türde ve 68 sınıfta toplam 1 360 000 görüntüden oluşmaktadır.

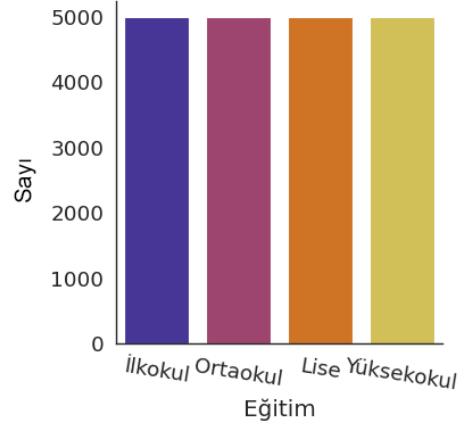
MFHD veri kümesi; Cinsiyet (bay: 10 000, bayan: 10 000), Eğitim durumu (ilkokul: 5 000, ortaokul: 5 000, lise: 5 000, yüksekokul: 5 000), Yaş grubu (5-11 yaş arası: 6 448, 12-19 yaş arası: 9 245, 20-30 yaş arası: 3 748, 31-65 yaş arası: 559) ve Hobi durumuna (kitap: 5 268, TV: 1 893, internet: 4 804, oyun: 1 945, spor: 2 285, müzik: 2 179, resim: 766, gezi: 860) göre 4 nitelik türünde sınıflandırılmıştır. Türkçe el yazısı karakter formlarının nitelik türlerinin dağılımları Şekil 3.1'de sunulmuştur.

a) Cinsiyet
Bay : 10 000
Bayan : 10 000

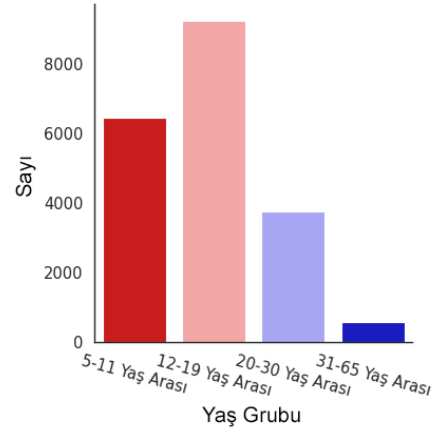


b) Eğitim Durumu

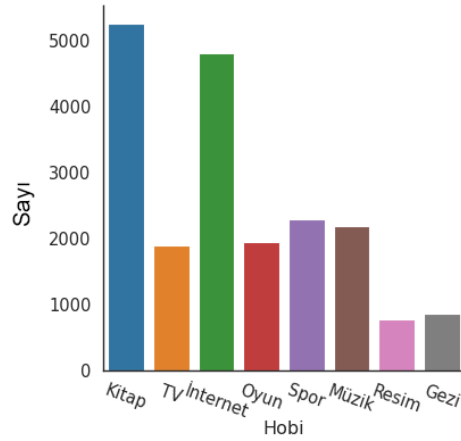
İlkokul	: 5 000
Ortaokul	: 5 000
Lise	: 5 000
Yüksekokul	: 5 000

**c) Yaş Grubu**

5-11 Yaş Arası	: 6 448
12-19 Yaş Arası	: 9 245
20-30 Yaş Arası	: 3 748
31-65 Yaş Arası	: 559

**d) Hobi Durumu**

Kitap	: 5 268
TV	: 1 893
İnternet	: 4 804
Oyun	: 1 945
Spor	: 2 285
Müzik	: 2 179
Resim	: 766
Gezi	: 860



Şekil 3.1 Formlarda bulunan nitelik sayılarının dağılımı, (a) cinsiyete göre (b) eğitim durumuna göre (c) yaş grubuna göre (d) hobi durumuna göre

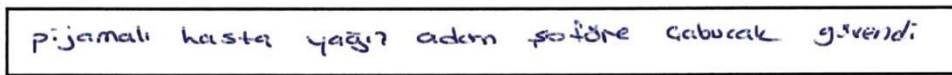
3.1. Materyal

3.1.1. MFHD Türkçe yazı şablonu hazırlama süreci

El yazısı karakter verileri; Türkçe Yazı Şablonu (TTT - Turkish Text Template) adını verdiğimiz, özel olarak tasarladığımız tek sayfalık A4 boyutunda özenle hazırlanmış formlar kullanılarak toplanmıştır. Formlar, Kişisel Verileri Koruma Kanunu (KVKK) (İnternet, 2019) dikkate alınarak gönüllülük esasına göre toplanmış ve kişisel veri içermemektedir. Gönüllü katılımcılardan, TTT formunda belirlenen karakter alanı sınırlarını aşmadan herhangi bir kalem türü ile yazmaları istenmiştir.

TTT formu; nitelik, küçük harf, büyük harf ve rakam olmak üzere 4 bölümden oluşmaktadır. Nitelik bölümünde yaş bilgisi yazmak için 1 alan bulunmakta olup diğer bölümleri işaretlemek için cinsiyet bilgisi 2 alan, eğitim durumu 8 alan ve hobi bilgisi 8 alandan oluşmaktadır. Küçük harf bölümünde, karakterler için 29 alan ve 1 pangram³ alanı, yine büyük harf bölümünde karakterler için 29 alan ve 1 pangram alanı, son olarak rakam bölümünde karakterler için 10 alan bulunmaktadır. TTT'de nitelik bilgisi ve karakter yazmak için toplam 89 alan yer alır. Pangram alanına yazılacak cümle, hiçbir anlam içermeyen rastgele kelimelerden oluşturulmuştur.

Küçük harf ve büyük harf pangramlar, yazı stilleri üzerine yapacağımız ileriki çalışmada kullanılacağından bu veri kümesinde paylaşılmamıştır. TTT formlarından rastgele seçilen küçük harf pangram alanı Şekil 3.2'de gösterilmiştir. TTT'nin rastgele seçilen, hizalanmamış örnek form görüntüsü en/boy oranı korunarak Şekil 3.3'de sunulmuştur.



Şekil 3.2 TTT formlarından rastgele seçilen küçük harf pangram alanı

³ Pangram, bir alfabede bulunan tüm harflerin kullanılmasıyla elde edilen kelime veya kelimeler bütünüdür

TÜRKÇE YAZI ŞABLONU

Yaş Cinsiyet Bay Bayan

Eğitim Durumu (X ile işaretleyiniz.)

Yok İlkokul Ortaokul Lise Ön Lisans Lisans Y. Lisans Doktora

Boş vakitlerinizi en çok ne ile değerlendirirsiniz? (Sadece 1(bir) sık işaretleyiniz.)

Kitap TV İnternet Oyun Spor Müzik Resim Gezi

ÖNEMLİ UYARI: HARFLER VE RAKAMLAR KENARLARA TEMAS ETMEDEN KUTU İÇİNE YAZILMALIDIR. SİLGİ KULLANILMADAN ve KARALAMA OLMADAN DİKKATLİCE YAZILMASI ÖNEMLE RİCA OLUNUR.

Küçük Harfler

a	b	c	ç	d	e	f	g	ğ	h
ı	İ	j	k	l	m	n	o	ö	p
r	s	ş	t	u	ü	v	y	z	

Cümle:* pijamalı hasta yağız adam şoföre çabucak güvendi.

Küçük Harf Yazın:

Büyük Harfler

A	B	C	Ç	D	E	F	G	Ğ	H
I	İ	J	K	L	M	N	O	Ö	P
R	S	Ş	T	U	Ü	V	Y	Z	

CÜMLE:* PİJAMALI HASTA YAĞIZ ADAM ŞOFÖRE ÇABUCAK GÜVENDİ.

BÜYÜK HARF YAZIN:

Rakamlar

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

Bilgi: Bu şablon Yüksek Lisans Tezinde DATASET (Veri Kümesi) olarak kullanılacaktır. Kişisel Verileri Koruma Kanunu dikkate alınarak hazırlanmıştır. * Yazılan Cümle 29 harfi içeren rastgele kelimelerden seçilmiştir.

Şekil 3.3 Doldurulmuş TTT formlarından rastgele seçilen bir örnek (hizalanmamış)

3.1.2. MFHD TTT toplama süreci

TTT formu kullanılarak; farklı seviyelerde eğitim öğretim faaliyeti yürüten özel ve resmi kuruluşlarda, üniversitelerde, sosyal faaliyet yürütülen alanlarda ve kamuya açık ortak alanlarda resmi izin ve etik kurulu kararı ile okuma yazma bilen her yaş ve cinsiyetten katılımcılardan bilgilendirme yapılarak gönüllülük esasına göre el yazısı karakter örnekleri toplanmıştır.

Veri toplama süreci; 14 farklı şehirde, 18 farklı üniversitede, ilkokul, ortaokul, lise seviyesinde resmi ve özel okullarda, kamuya açık alanlarda yaklaşık 24 ay süreyle gönüllü araştırmacı üniversite öğrencileri, gönüllü öğretmenler ve yazarlar tarafından bizzat bilgilendirme yapılarak gönüllülük esasına göre gözetim altında toplanmıştır. Gönüllü katılımcılardan, TTT formunu doldururken karalama yapılmaması ve silgi kullanılmaması istenmiştir. TTT formları, gönüllü ve istekli katılımcılar tarafından ortalama 7 dakikada doldurulmuştur.

3.1.3. MFHD veri kümesi hazırlık süreci

MFHD veri kümesi hazırlama sürecinde 29 500 tekil katılımcıdan TTT kullanılarak el yazısı karakter formları toplanmıştır. Elle yapılan ilk incelemelerde yetersiz ve eksik doldurulan, çok fazla gürültü ve karalama içeren formlar elenmiştir. 25 150'ye düşürülen bu formlar kontrollü bir şekilde profesyonel tarayıcıda 300 dpi çözünürlükte RGB olarak taranıp dijital ortama aktarılmıştır. Taramadan kaynaklanan gürültülü ve düzensiz formlar da elendikten sonra bu sayı 21 300 adete inmiştir. Toplanan el yazısı formlar, dijital ortamda türüne göre ayrıştırıldıktan sonra tüm nitelik (cinsiyet, yaş, eğitim ve hobi) sayılarında eşit dağılım oluşmadığı tespit edilmiştir. Cinsiyet ve eğitim durumu nitelik sayılarının eşit dağılım oluşturması için 20 000 adetinin nihai veri kümesi olarak kullanılması kararlaştırılmıştır. 20 000 Türkçe el yazısı karakter formununun 4 farklı nitelik türüne göre dağılımı Şekil 3.1'de gösterilmiştir.

3.2. Yöntem

MFHD veri kümesini oluşturmak için belirlenmiş dijital ortamda bulunan 20 000 el yazısı karakter formu, aşağıda belirtilen adımlar uygulanarak işlenmiş, etiketlenmiş ve araştırmacıların kullanabileceği bir formata dönüştürülmüştür.

1. Form nitelik bilgisi oluşturma
2. Form hizalama
3. Form gürültü temizleme
4. Form karakter algılama, çıkarma ve RGB gri-ton dönüşümü
5. Renkleri ters çevirme
6. Karakter boyut normalleştirme
7. Karakter etiketleme

3.2.1. Form nitelik bilgisi oluşturma

Bu adımda, dijitalleştirilen 20 000 el yazısı karakter formunun nitelik bilgisi oluşturulmuştur. Her form, hazırladığımız özel yazılım ve dijital ortamda fiziksel olarak gözle incelenip nitelik bilgileri çıkarılmıştır. Türkçe el yazısı karakter formlarında bulunan cinsiyet, yaş, eğitim ve hobi gibi nitelik bilgileri MFHD-Feature isimli CSV dosyasında tutulmuştur. Nitelik dosyası 20 001 satır ve 5 sütundan oluşmaktadır. İlk satırda etiket bilgileri yer almakta olup takip eden satırlarda her formun nitelik bilgileri yer almaktadır. MFHD nitelik dosyasının içeriği Şekil 3.4'te gösterilmiştir.

```
MFHD-Feature.csv
Form-id,Yaş,Cinsiyet,Eğitim,Hobi
1,8,Bay,İlkokul,Spor
2,8,Bay,İlkokul,Spor
.
20000,28,Bayan,Yüksekokul,TV
```

Şekil 3.4 MFHD-Feature dosya içeriği

Türkçe el yazısı karakter formlarının nitelik bilgilerinin saklandığı MFHD-Feature.csv dosyasında bulunan etiket bilgileri aşağıda açıklanmıştır.

Form-id: 1'den 20 000'e kadar olan sayılardan oluşmaktadır. Formu dolduran katılımcı bilgisini ifade eder.

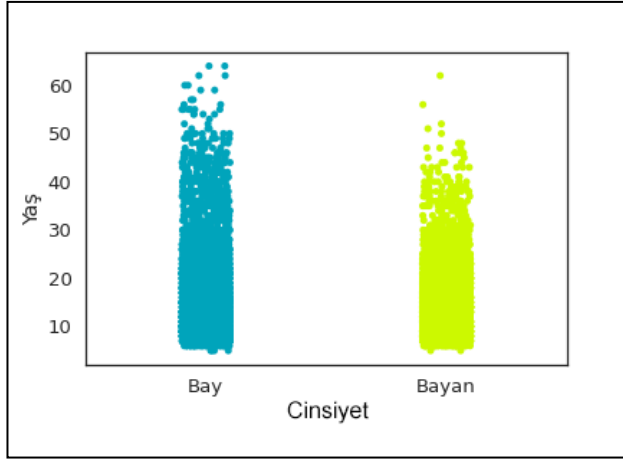
Yaş: 5 ile 64 arasında değişen sayılardan oluşmaktadır. Formu dolduran katılımcının yaş bilgisini ifade eder. Katılımcıların niteliklere göre yaş dağılımı Şekil 3.5'te gösterilmiştir.

Cinsiyet: Bay ve bayan bilgisinden oluşmaktadır. Formu dolduran katılımcının cinsiyet bilgisini ifade eder. Katılımcılar cinsiyete göre; 10 000 bay ve 10 000 bayan şeklinde dağılım göstermektedir.

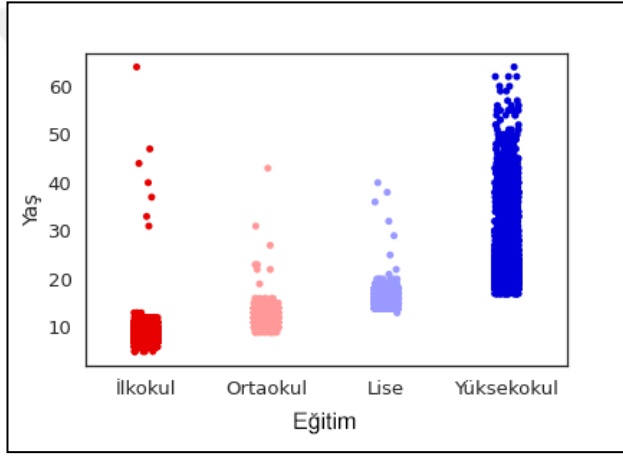
Eğitim: İlkokul, ortaokul, lise ve yüksekokul bilgilerinden oluşmaktadır. Formu dolduran katılımcının eğitim durumunu ifade eder. Eğitim durumunun eşit şekilde dağılımının sağlanması için ön lisans, lisans, yüksek lisans ve doktora seviyesindeki katılımcılar "yüksekokul" olarak etiketlenmiştir. Katılımcılar eğitim durumuna göre; 5 000 ilkök, 5 000 ortaokul, 5 000 lise ve 5 000 yüksekokul şeklinde dağılım göstermektedir.

Hobi: Kitap, TV, internet, oyun, spor, müzik, resim ve gezi bilgilerinden oluşmaktadır. Formu dolduran katılımcının hobi bilgisini ifade eder. Katılımcılar hobi bilgisine göre; 5 268 kitap, 1 893 TV, 4 804 internet, 1 945 oyun, 2 285 spor, 2 179 müzik, 766 resim ve 860 gezi şeklinde dağılım göstermektedir.

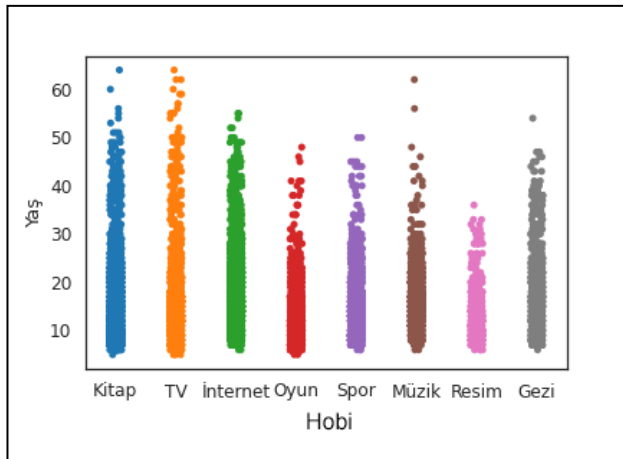
Erikson'un psikososyal gelişim kuramına göre bireyin gelişim dönemleri de dikkate alınarak veri kümesindeki yaş bilgisi, deneyler için 4 grupta (5-11 yaş arası, 12-19 yaş arası, 20-30 yaş arası ve 31-65 yaş arası) toplanmıştır (Erikson, 1993). Katılımcılar yaş grubuna göre; 6 448 kişi 5-11 yaş arası, 9 245 kişi 12-19 yaş arası, 3 748 kişi 20-30 yaş arası ve 559 kişi 31-65 yaş aralıklarında dağılım göstermektedir. MFHD veri kümesinin cinsiyet, eğitim durumu ve yaş grubuna göre dağılım grafiği Şekil 3.6'te gösterilmiştir.



(a)

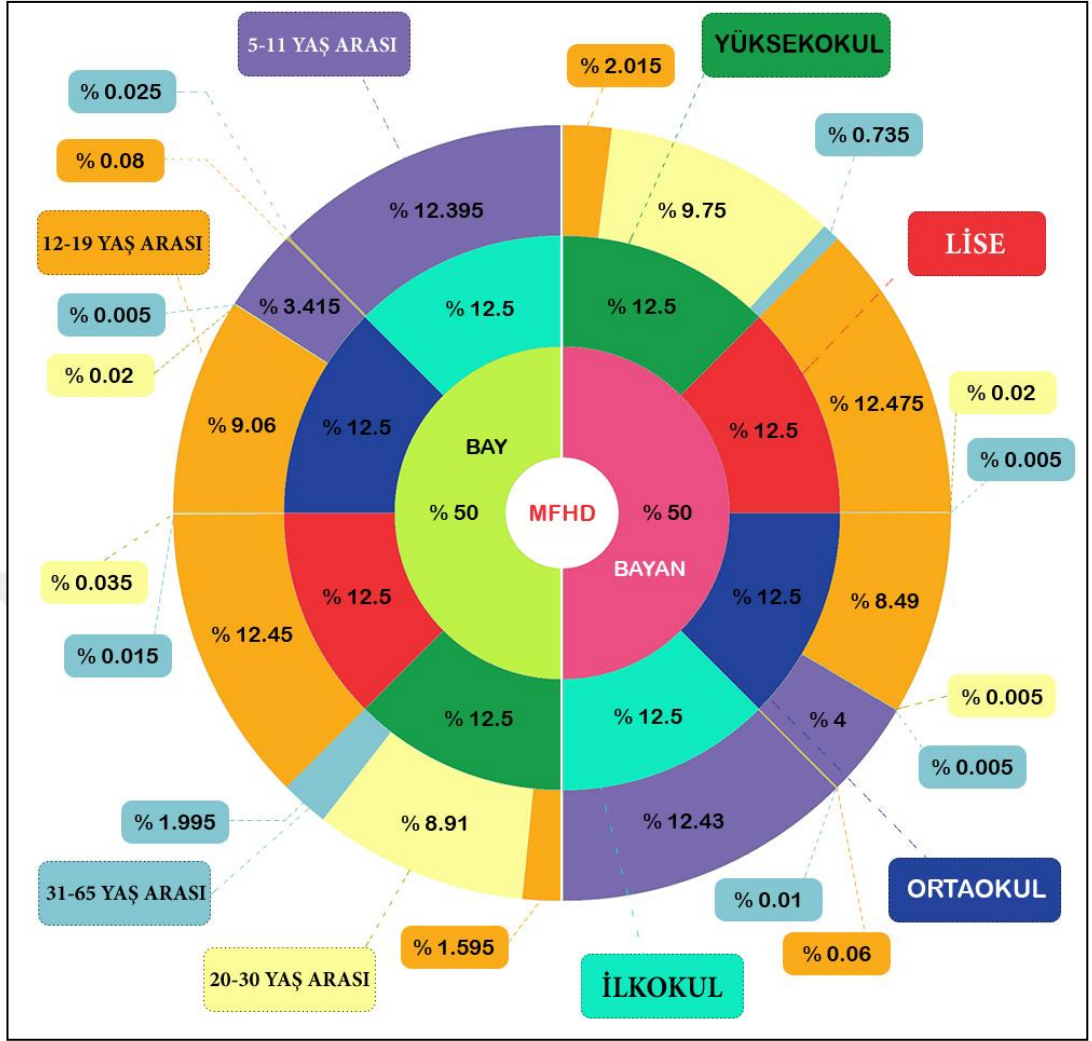


(b)



(c)

Şekil 3.5 Katılımcıların nitelik yaş dağılımları. (a) cinsiyete göre (b) eğitim durumuna göre (c) hobi durumuna göre



Şekil 3.6 MFHD veri kümesinin cinsiyet, eğitim durumu ve yaş grubuna göre dağılım grafiği

Deneyler sırasında doğru sınıflandırma yapılabilmesi için (yaş grubu ve hobi durumunda eşit dağılım (Bias) durumunu sağlayabilmek için) 20 000 form içerisinde yaş grubu ve hobi durumunda en küçük nitelik sayısına göre rastgele örneklem alınmıştır. Bu örneklem veri kümeleri yaş grubu için; küçük harf, büyük harf ve rakamlara göre sırasıyla MFHD-LA, MFHD-UA ve MFHD-DA olarak adlandırılmıştır. Aynı şekilde hobi için kullanılan örneklem veri kümeleri de küçük harf, büyük harf ve rakamlara göre sırasıyla MFHD-LH, MFHD-UH ve MFHD-DH olarak adlandırılmıştır. Yaş grubu ve hobi durumuna göre hazırlanan bu örneklem veri kümeleri de ayrıca paylaşılmıştır.

3.2.2. Form hizalama

TTT ile toplanan formlar tarayıcı ile dijital ortama aktarılırken bazılarında gözle görülebilen bazılarında ise gözle görülemeyen eksen kaymaları (yatay ve dikey yönde) olduğu yapılan incelemelerde tespit edilmiştir. El yazısı karakterlerinin doğru alanlardan eksiksiz olarak alınabilmesi ve analizinin hatasız yapılabilmesi için formların hizalanması gerektiği görülmüştür.

Hizalama için Python kütüphanelerinde bulunan OpenCV içindeki Homografi metodu kullanılmıştır. Homografi, bir görüntüdeki noktaları başka bir görüntüde karşılık gelen noktaya eşleyen bir matristir (Prathap ve ark., 2016). TTT el yazısı karakter formunda dikkat çeken en belirgin özellikler çizgiler ve köşelerdir. İki görüntü arasında iyi bir özellik eşleştirmesi sağlamak için köşeler karşılaştırılabilir. Köşeler, görüş açılarındaki değişikliklere göre daha kararlı özellikler içermesi ve komşuluğunda ani bir yoğunluk değişikliği göstermesinden dolayı görüntülerin eşleştirilebilecek en iyi özelliklerinden birisidir (Vaghela ve Naina, 2014). Çeşitli köşe algılama algoritmaları kullanılarak görüntülerde köşeler tespit edilebilir. Köşe algılama algoritmalarından bazıları, Harris algoritması (Harris ve Stephens, 1988), SUSAN algoritması (Smith ve Brady, 1997), makine öğrenmesi tabanlı FAST algoritması (Trajković ve Hedley, 1998), SIFT (Scale Invariant Feature Transform) algoritması (Lowe, 2004), SURF (Speeded Up Robust Feature) algoritması (Bay ve ark., 2006) ve ORB algoritmasıdır (Rublee ve ark., 2011).

TTT formlarının hizalanması için yaptığımız denemeler ve literatürde yapılan çalışmalar da dikkate alınarak özellik çıkarma ve köşeleri belirleme için ORB algoritması kullanılmıştır. Python ve kütüphanelerini kullanarak hazırladığımız yazılımla, orijinal TTT formu ile taranıp dijital ortama aktarılan tüm el yazısı karakter formları karşılaştırıldıktan sonra eksen kayması bulunan formlar belirlenmiş ve hizalanmıştır. Şekil 3.7'da eksen kayması bulunan rastgele bir form, Şekil 3.8'de ise hizalanmış aynı form gösterilmektedir. Şekil 3.8'deki siyah şeritler formun hangi yönde ve ne kadar hizalandığını göstermektedir.

TÜRKÇE YAZI ŞABLONU

Yas **Cinsiyet** Bay Bayan

Eğitim Durumu (X ile işaretleyiniz.)

Yok İlkokul Ortaokul Lise Ön Lisans Lisans Y. Lisans Doktora

Bos vakitlerinizi en çok ne ile değerlendirirsiniz? (Sadece 1(bir) sık işaretleyiniz.)

Kıtap TV İnternet Oyun Spor Müzik Resim Gezi

ÖNEMLİ UYARI: HARFLER VE RAKAMLAR KENARLARA TEMAS ETMEDEN KUTU İÇİNE YAZILMALIDIR. SİLGİ KULLANILMADAN ve KARALAMA OLMADAN DİKKATLİCE YAZILMASI ÖNEMLE RICA OLUNUR.

Küçük Harfler

a	b	c	ç	d	e	f	g	ğ	h
ı	î	ı	k	l	m	n	o	ö	p
r	s	ş	t	u	ü	v	y	z	

Cümle:* pijamalı hasta yağız adam şoföre çabucak güvendi.

Küçük Harf Yazın: pijamalı hasta yağız adam şoföre çabucak güvendi.

Büyük Harfler

A	B	C	Ç	D	E	F	G	Ğ	H
I	İ	J	K	L	M	N	O	Ö	P
R	S	Ş	T	U	Ü	V	Y	Z	


CÜMLE:* PİJAMALI HASTA YAĞIZ ADAM ŞOFÖRE ÇABUCAK GÜVENDİ.

BÜYÜK HARF YAZIN: PİJAMALI HASTA YAĞIZ ADAM ŞOFÖRE ÇABUCAK GÜVENDİ.

Rakamlar

0	1	2	3	4	5	6	7	8	9
---	---	---	---	---	---	---	---	---	---

Bilgi: Bu şablon Yüksek Lisans Tezinde DATASET (Veri Kümesi) olarak kullanılacaktır. Kişisel Verileri Koruma Kanunu dikkate alınarak hazırlanmıştır. * Yazılan Cümle 29 harfi içeren rastgele kelimelerden seçilmiştir.

DÖNDÜRME 

Şekil 3.7 Eksen kayması bulunan rastgele örnek bir form (En/boy oranı korunmuş)

TÜRKÇE YAZI ŞABLONU

Yaş Cinsiyet Bay Bayan

Eğitim Durumu (X ile işaretleyiniz.)
 Yok İlkokul Ortaokul Lise Ön Lisans Lisans Y. Lisans Doktora

Bos vakitlerinizi en çok ne ile değerlendirirsiniz? (Sadece 1(bir) sık işaretleyiniz.)
 Kitap TV İnternet Oyun Spor Müzik Resim Gezi

ÖNEMLİ UYARI: HARFLER VE RAKAMLAR KENARLARA TEMAS ETMEDEN KUTU İÇİNE YAZILMALIDIR. SİLGİ KULLANILMADAN ve KARALAMA OLMADAN DİKKATLİCE YAZILMASI ÖNEMLE RİCA OLUNUR.

Küçük Harfler

a	b	c	ç	d	e	f	g	ğ	h
ı	î	j	k	l	m	n	o	ö	p
r	s	ş	t	u	ü	v	y	z	

Cümle:* pijamalı hasta yağız adam şoföre çabucak güvendi.

Küçük Harf Yazın: pijamalı hasta yağız adam şoföre çabucak güvendi.

Büyük Harfler

A	B	C	Ç	D	E	F	G	Ğ	H
I	İ	J	K	L	M	N	O	Ö	P
R	S	Ş	T	U	Ü	V	Y	Z	

CÜMLE:* PİJAMALI HASTA YAĞIZ ADAM ŞOFÖRE ÇABUCAK GÜVENDİ.

BÜYÜK HARF YAZIN: PİJAMALI HASTA YAĞIZ ADAM ŞOFÖRE ÇABUCAK GÜVENDİ

Rakamlar

0	1	2	3	4	5	6	7	8	9
---	---	---	---	---	---	---	---	---	---

Bilgi: Bu şablon Yüksek Lisans Tezinde DATASET (Veri Kümesi) olarak kullanılacaktır. Kişisel Verileri Koruma Kanunu dikkate alınarak hazırlanmıştır. * Yazılan Cümle 29 harfi içeren rastgele kelimelerden seçilmiştir.

Şekil 3.8 Eksen kayması olan hizalanmış aynı form (Bakınız Şekil 3.7)

3.2.3. Form gürültü temizleme

TTT ile toplanan formlar, katılımcıların formu doldururken oluşturduğu gürültü ile birlikte profesyonel tarayıcılarla dijital ortama aktarılırken bazen tarayıcı kaynaklı bazen de kâğıtlar arasında bulunan toz zerreciklerinden dolayı gürültü içerebilmiştir. Çeşitli gürültü temizleme ve giderme algoritmaları kullanılarak görüntülerde bulunan gürültüler giderilebilir. Literatürde bilinen gürültü temizleme ve giderme algoritmalarından bazıları, Gauss filtresi (Gaussian filter), Ortalama filtre (Mean filter), Medyan filtre (Median filter) ve İkili filtre (Bilateral filter)'dir.

TTT formlarının varsa gürültülerinin giderilmesi veya azaltılması için yaptığımız denemeler ve literatürde yapılan çalışmalar da dikkate alınarak Python kütüphanelerinde bulunan OpenCV içindeki Medyan filtre (Median filter, 3x3 kernel) metodu kullanılmıştır. Medyan filtreleme, gürültü bastırma için yararlı olan doğrusal olmayan bir görüntü işleme tekniğidir (Justusson, 1981). Medyan filtrelemede her piksel, bir dizi komşu pikselin değerinin medyanı ile değiştirildiğinden dolayı bu yöntem, filtre doğrusal olmasa da bir filtreleme tekniği olarak kabul edilebilir (Saxena ve Kourav, 2014). Medyan filtresi basit ve sıra istatistiklerine dayanan güçlü ve doğrusal olmayan bir filtre olmakla birlikte bir piksel ile diğer piksel arasındaki yoğunluk değişimi miktarını azaltmak için de kullanılır (Hambal ve ark., 2017).

Python ve kütüphanelerini kullanarak hazırladığımız yazılım ile tüm formlar gürültü temizleme işlemine tabi tutulmuş ve varsa form üzerinde bulunan gürültüler büyük çoğunlukta giderilmiştir. Şekil 3.9'de gürültülü rastgele bir form ve Şekil 3.10'de gürültüsü giderilmiş aynı form gösterilmektedir.

TÜRKÇE YAZI ŞABLONU

Yaş Cinsiyet Bay Bayan

Eğitim Durumu (X ile işaretleyiniz.)

Yok İlkokul Ortaokul Lise Ön Lisans Lisans Y. Lisans Doktora

Boş vakitlerinizi en çok ne ile değerlendirirsiniz? (Sadece 1(bir) sık işaretleyiniz.)

Kitap TV İnternet Oyun Spor Müzik Resim Gezi

ÖNEMLİ UYARI: HARFLER VE RAKAMLAR KENARLARA TEMAS ETMEDEN KUTU İÇİNE YAZILMALIDIR. SİLĞİ KULLANILMADAN ve KARALAMA OLMADAN DİKKATLİCE YAZILMASI ÖNEMLİ RİCA OLUNUR.

Küçük Harfler

a	b	c	ç	d	e	f	g	ğ	h
a	b	c	ç	d	e	f	g	ğ	h
i	ı	j	k	l	m	n	o	ö	p
i	ı	j	k	l	m	n	o	ö	p
r	s	ş	t	u	ü	v	y	z	
r	s	ş	t	u	ü	v	y	z	

Cümle:* pijamalı hasta yağız adam şoföre çabucak güvendi.

Küçük Harf Yazın: pijamalı hasta yağız adam çabucak güvendi.

Büyük Harfler

A	B	C	Ç	D	E	F	G	Ğ	H
A	B	C	Ç	D	E	F	G	Ğ	H
I	İ	J	K	L	M	N	O	Ö	P
I	İ	J	K	L	M	N	O	Ö	P
R	S	Ş	T	U	Ü	V	Y	Z	
R	S	Ş	T	U	Ü	V	Y	Z	

CÜMLE:* PİJAMALI HASTA YAĞIZ ADAM ŞOFÖRE ÇABUCAK GÜVENDİ.

BÜYÜK HARF YAZIN: PİJAMALI HASTA YAĞIZ ADAM ŞOFÖRE ÇABUCAK GÜVENDİ

Rakamlar

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

Bilgi: Bu şablon Yüksek Lisans Tezinde DATASET (Veri Kümesi) olarak kullanılacaktır. Kişisel Verileri Koruma Kanunu dikkate alınarak hazırlanmıştır. * Yazılan Cümle 29 harfi içeren rastgele kelimelerden seçilmiştir.

Şekil 3.9 Rastgele seçilen gürültülü bir form (hizalanmamış)

TÜRKÇE YAZI ŞABLONU

Yaş Cinsiyet Bay Bayan

Eğitim Durumu (X ile işaretleyiniz.)
 Yok İlkokul Ortaokul Lise Ön Lisans Lisans Y. Lisans Doktora

Bos vakitlerinizi en çok ne ile değerlendirirsiniz? (Sadece 1(bir) şık işaretleyiniz.)
 Kitap TV İnternet Oyun Spor Müzik Resim Gezi

ÖNEMLİ UYARI: HARFLER VE RAKAMLAR KENARLARA TEMAS ETMEDEN KUTU İÇİNE YAZILMALIDIR. SİLGİ KULLANILMADAN ve KARALAMA OLMADAN DİKKATLİCE YAZILMASI ÖNEMLE RİCA OLUNUR.

Küçük Harfler

a	b	c	ç	d	e	f	g	ğ	h
ā	b	c	ç	d	e	f	g	ğ	h
i	ı	j	k	l	m	n	o	ö	p
ī	ı	j	k	l	m	n	o	ö	p
r	s	ş	t	u	ü	v	y	z	
r	s	ş	t	u	ü	v	y	z	

Cümle:* pijamalı hasta yağız adam şoföre çabucak güvendi.

Küçük Harf Yazın: pijamalı hasta yağız adam çabucak güvendi.

Büyük Harfler

A	B	C	Ç	D	E	F	G	Ğ	H
A	B	C	Ç	D	E	F	G	Ğ	H
I	İ	J	K	L	M	N	O	Ö	P
I	İ	J	K	L	M	N	O	Ö	P
R	S	Ş	T	U	Ü	V	Y	Z	
R	S	Ş	T	U	Ü	V	Y	Z	

CÜMLE:* PİJAMALI HASTA YAĞIZ ADAM ŞOFÖRE ÇABUCAK GÜVENDİ.

BÜYÜK HARF YAZIN: PİJAMALI HASTA YAĞIZ ADAM ŞOFÖRE ÇABUCAK GÜVENDİ

Rakamlar

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

Bilgi: Bu şablon Yüksek Lisans Tezinde DATASET (Veri Kümesi) olarak kullanılacaktır. Kişisel Verileri Koruma Kanunu dikkate alınarak hazırlanmıştır. * Yazılan Cümle 29 harfi içeren rastgele kelimelerden seçilmiştir.

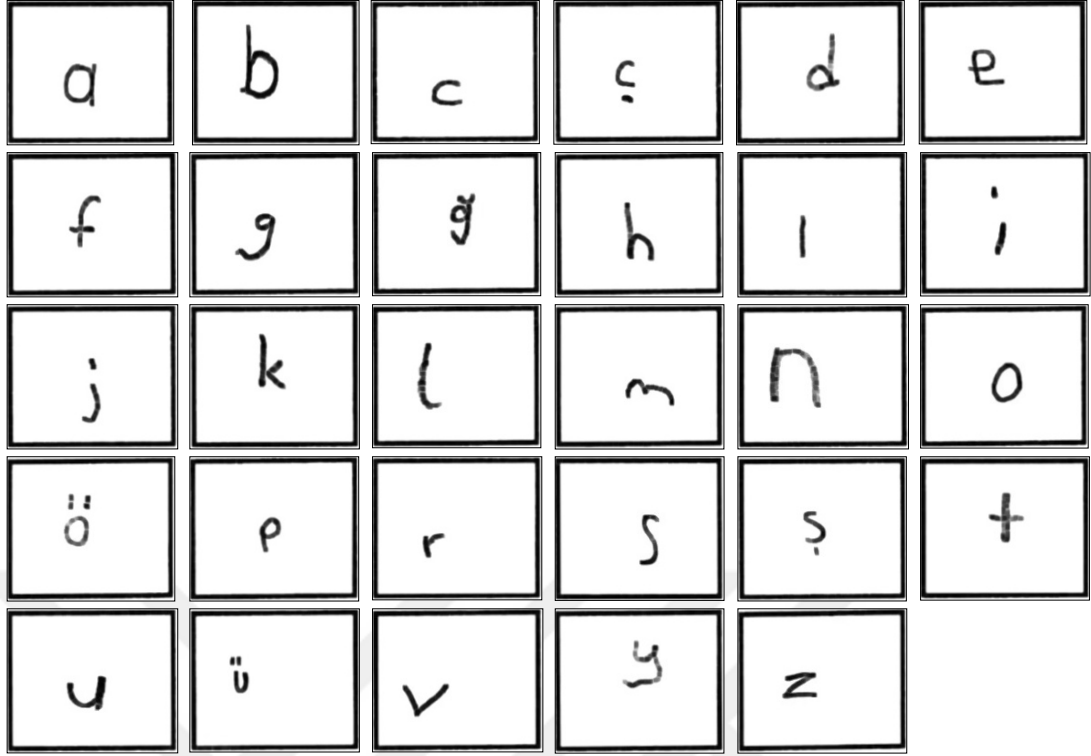
Şekil 3.10 Gürültüsü giderilmiş aynı (Bakınız Şekil 3.9) form (hizalanmış)

3.2.4. Form karakter algılama, çıkarma ve RGB gri-ton dönüşümü

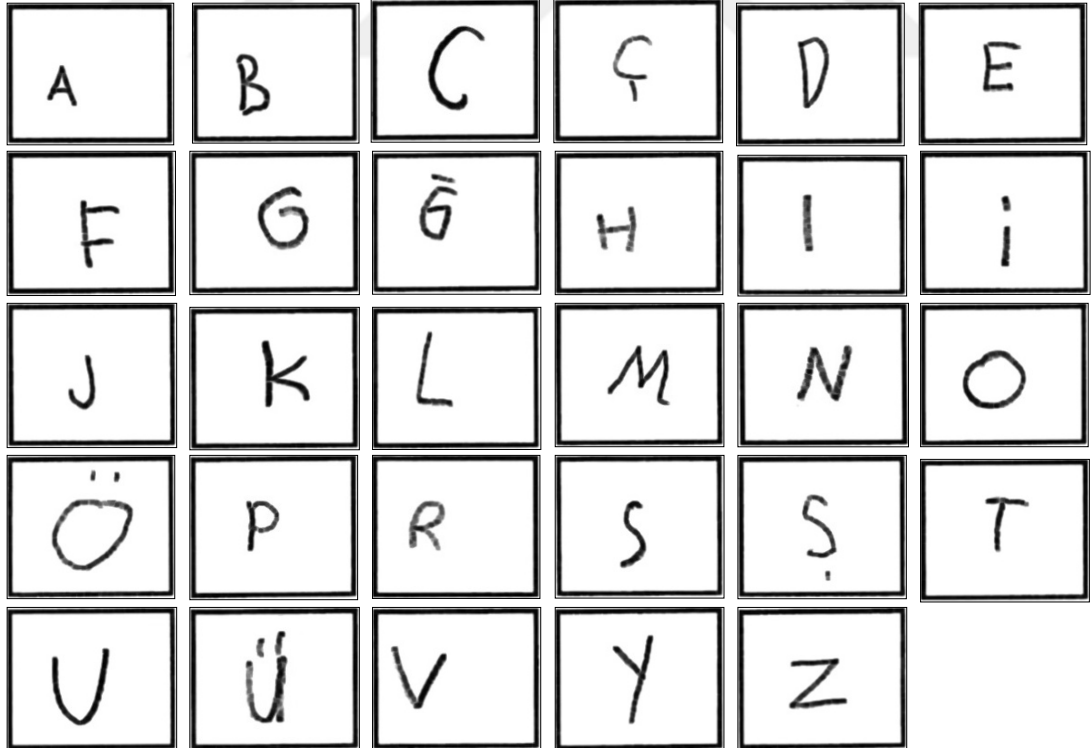
Tarayıcı ile taranıp dijital ortama aktarılan formlar hizalandıktan ve varsa gürültüleri giderildikten sonra morfolojik yöntemlerle karakterler çıkartılmıştır. TTT formunda, el yazısı karakterlerin çevresinde bulunan kenarlar kullanılarak karakterler çıkarılmıştır.

Kenarlar; parlaklığı, renkteki süreksizlikleri ve iki yüzeyin kesişimini temsil eder (Hansen ve Gegenfurtner, 2017). Kenar algılamasında önemli aşamalar kaydedilmiştir ancak insan görmesi ile karşılaştırıldığında hala eksiklikler bulunmaktadır (Yang ve ark., 2022). Literatürde kenar ve nesne kontur tespiti için Roberts operatörü (Roberts, 1963), Sobel operatörü (Sobel ve Feldman, 1968) ve Gabor enerji filtreleri (Boukerroui ve ark., 2004) gibi yöntemler içeren bazı öncü çalışmalar yapılmıştır. Bunun yanında görüntü bölümlenme için çok sayıda yöntem (Peng ve ark., 2013; Kaur ve Kaur, 2014; Khan, 2014; Minaee ve ark., 2021) önerilmiştir. Bunlar temel olarak kümeleme tabanlı, havza tabanlı, eşik tabanlı, çizge teorisi tabanlı ve derin öğrenme tabanlı yöntemler olarak ayrılırlar. Bilgisayarlı görmenin temel görevlerinden biri olan görüntü segmentasyonuna kapsamlı bir giriş çalışmamız kapsamı dışında olduğundan ayrıntılı bilgi için ilgili kaynaklardan yararlanılabilir.

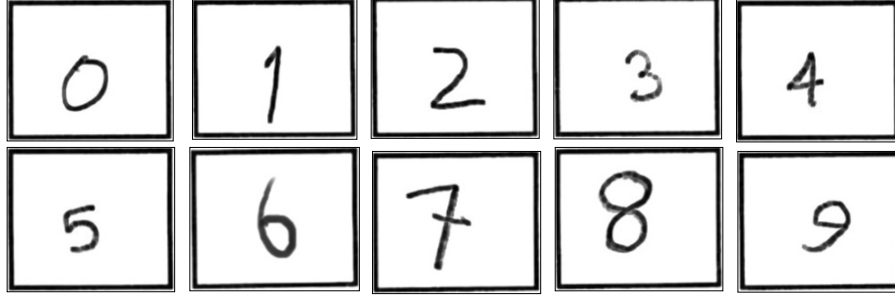
Python kütüphanelerinde bulunan OpenCV kullanarak hazırladığımız yazılım ile yatay ve dikey eksende bulunan karakter alanları belirlenmiş ve bunlara bir kod numarası verilmiştir. Verilen kod numarasına karşılık gelen küçük harf (29), büyük harf (29) ve rakam (10) olmak üzere toplam 68 karakter alanı çıkarılmış ve numaralandırılmıştır. Çıkarılan karakterler RGB'den gri-ton'a (gray scale) dönüştürülerek saklanmıştır. Formlardan çıkarılan rastgele küçük harf (29 sınıf) karakter örnekleri Şekil 3.11'da, büyük harf (29 sınıf) karakter örnekleri Şekil 3.12'de ve rakam (10 sınıf) karakter örnekleri Şekil 3.13'de gösterilmiştir.



Şekil 3.11 Formlardan çıkarılan rastgele küçük harf (29 sınıf) karakter örnekleri



Şekil 3.12 Formlardan çıkarılan rastgele büyük harf (29 sınıf) karakter örnekleri

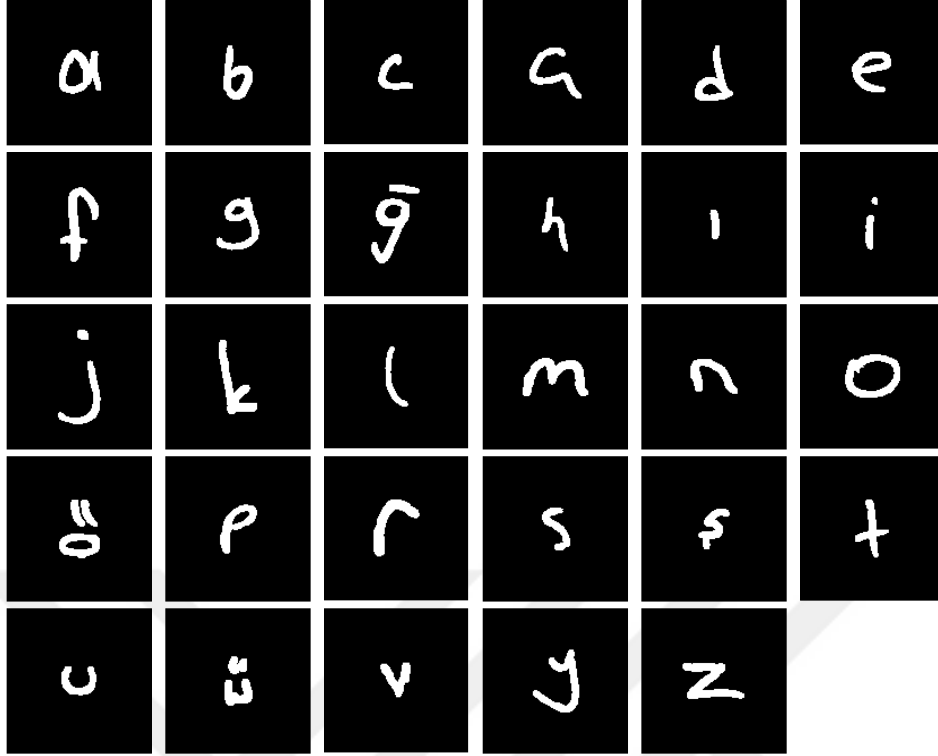


Şekil 3.13 Formlardan çıkarılan rastgele rakam (10 sınıf) karakter örnekleri

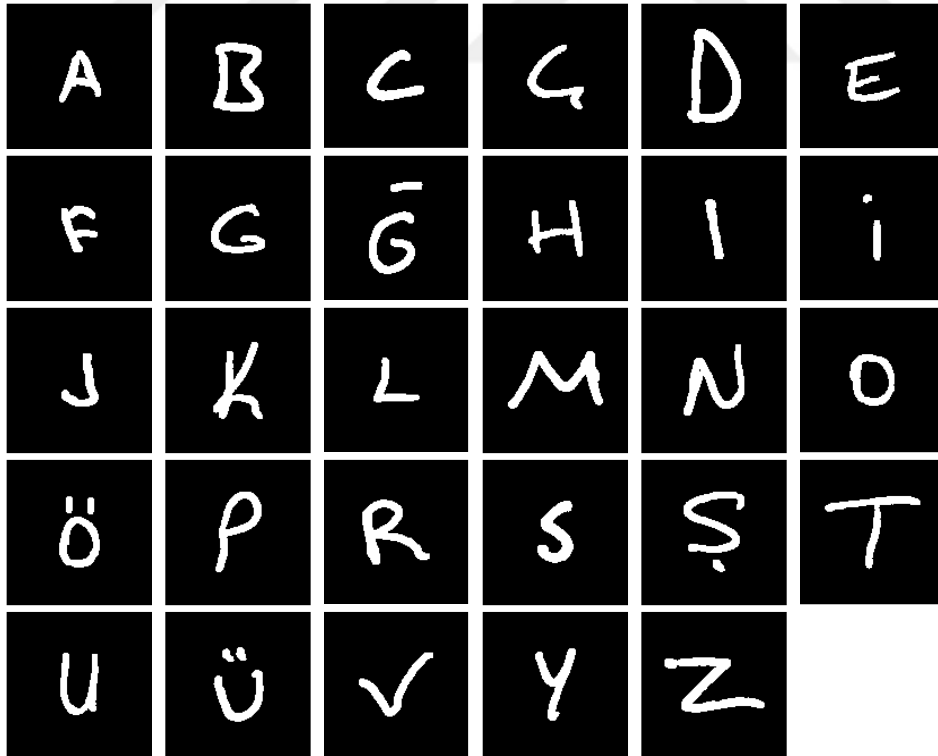
3.2.5. Renkleri ters çevirme

Çıkarılan karakterlerin çevresinde bulunan siyah kenarlıklar (Bakınız Şekil 3.11, Şekil 3.12 ve Şekil 3.13), Python ve kütüphanelerini kullanarak hazırladığımız yazılım ile kaldırıldıktan sonra renkleri eşikleme ile ters çevrilmiştir. Eşikleme, görüntüden alınan piksel değerinin belirlenen eşik değeri ile karşılaştırılarak yeni piksel değerlerinin atanması olan bir tekniktir. Eşik, her iki yönlü, yani eşikğin altında veya üstünde iki bölgede de bulunan bir değerdir. Eşiklemenin amacı, ön plan olarak kabul edilen bir nesneyi arka planından ayırmaktır (Pare ve ark., 2020). Eşik tabanlı bölütleme; literatürde bulunan çeşitli bölütleme teknikleri arasında, daha az hesaplama maliyeti ve daha yüksek verimliliğe sahip olduğu için en basit görüntü bölütleme yaklaşımlarından biridir (Sezgin ve Sankur, 2004)

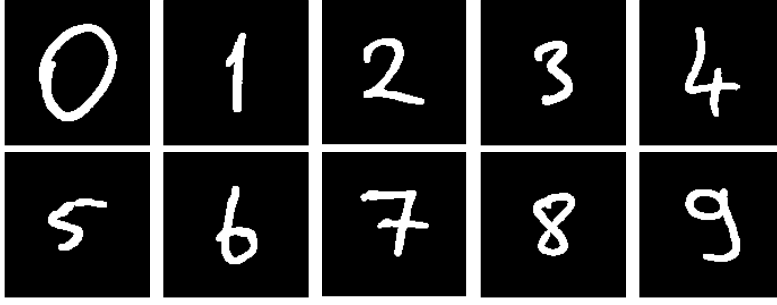
TTT formlarından çıkarılan karakterlerin çevresinde bulunan siyah kenarlıklar kaldırıldıktan sonra eşikleme için yaptığımız denemeler ve literatürde yapılan çalışmalar da dikkate alınarak Python kütüphanelerinde bulunan OpenCV kullanarak hazırladığımız yazılım ile renkler ters çevrilmiştir. Renkleri ters çevrilen karakterler çerçeve içine ortalanmış, derinliği 8 bit'e düşürülmüş ve .png formatında (128x128 piksel boyutlarında) saklanmıştır. Renkleri ters çevrilen rastgele küçük harf (29 sınıf) karakter örnekleri Şekil 3.14'de, büyük harf (29 sınıf) karakter örnekleri Şekil 3.15'de ve rakam (10 sınıf) karakter örnekleri Şekil 3.16'te gösterilmiştir.



Şekil 3.14 Renkleri ters çevrilen rastgele küçük harf (29 sınıf) karakter örnekleri



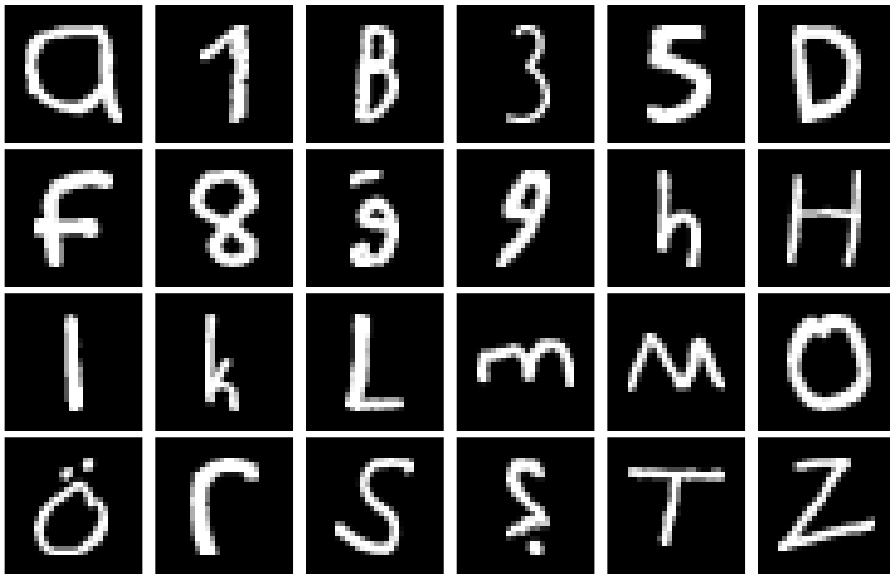
Şekil 3.15 Renkleri ters çevrilen rastgele büyük harf (29 sınıf) karakter örnekleri



Şekil 3.16 Renkleri ters çevrilen rastgele rakam (10 sınıf) karakter örnekleri

3.2.6. Karakter boyut normalleştirme

Siyah zemin üzerinden (Bakınız Şekil 3.14, Şekil 3.15 ve Şekil 3.16) beyaz renkli karakterlerin vektörel kontur hesaplaması yapılmış ve gerçek en-boy oranları korunarak 20x20 piksel boyutlarında normalleştirilmiştir (Lecun ve ark., 1998). Normalleştirme algoritmasının çıktıları gri renk seviyelerinden oluşmakta olup 20x20 piksele normalleştirilen karakterlerin ağırlık merkezleri 28x28 boyutundaki alan merkezlerine yerleştirilmesi yapılmıştır. Böylelikle MFHD veri kümesindeki görüntüler sık kullanılan (Benchmark Dataset) MNIST ile aynı özelliklere (Görüntü boyutu, gri-ton değeri) sahip olmaktadır. Boyutu normalleştirilen rastgele karakter (küçük harf, büyük harf ve rakam) örnekleri Şekil 3.17’de gösterilmiştir.



Şekil 3.17 Boyutu normalleştirilen rastgele karakter örnekleri (28x28 piksel boyutunda)

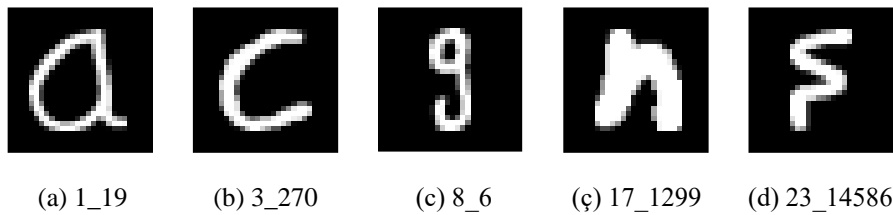
Çizelge 2.1’de verilmiş olan bazı veri kümelerinde çeşitli döndürme ve hizalama işlemi yapılırken MFHD veri kümesinde bilgi çıkarılma amaçlandığı için katılımcıların el yazı bilgisini bozmamak amacıyla TTT formlarından elde edilen karakterlere herhangi bir döndürme veya hizalama işlemi yapılmamıştır.

3.2.7. Karakter etiketleme

Çok büyük sayıda görüntü veya dosyanın bulunduğu veri kümelerinde kullanılabilirlik açısından etiketleme önemli bir işlem adımı olarak görülmektedir. Etiket düzeni, iletilen bilgilerin verimli ve doğru anlaşılmasında çok önemli bir rol oynamaktadır (Cmolik ve ark., 2020). MFHD veri kümesinin araştırmacılar tarafından kolay bir şekilde kullanılabilmesi, anlaşılabilir ve net olması için her görüntü ilgili veri kümesine göre etiketlenmiştir.

MFHD veri kümesi, Türk alfabesinde bulunan Latin harflerden oluşmaktadır. Veri kümesi 3 farklı türde (küçük harf, büyük harf ve rakam) sınıflandırılmış görüntü kümelerinden ve 1 nitelik dosyasından oluşmaktadır. Her veri kümesi kendi içinde ayrı ayrı etiketlenmiştir. Bu sınıflandırmalar aşağıda belirtilmiş olup çalışmanın bundan sonraki kısımlarında bu isimler kullanılacaktır.

MFHD-L (MFHD-Lower Case): MFHD veri kümesinde bulunan küçük harfleri içermektedir. Türk alfabesinde bulunan 'a' dan 'z' ye küçük harflerden oluşmaktadır. 1’den 29’a kadar olan sayılarla, Türk alfabesinde bulunan 29 küçük harf sırasına göre eşleştirilmiştir. Örneğin; a harfi '1' ile z harfi '29' ile etiketlenmiştir. MFHD-L veri kümesine ait etiketlenmiş rastgele karakter örnekleri Şekil 3.18’de gösterilmiştir.

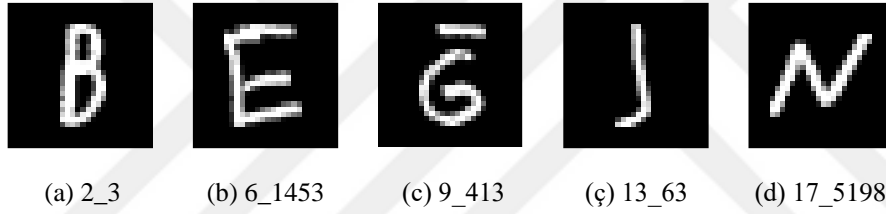


Şekil 3.18 Etiketlenmiş rastgele küçük harf karakter örnekleri.

Şekil 3.18 Açıklaması:

- (a): 1 harf sıra bilgisini, 19 Form-id bilgisini ifade etmektedir,
- (b): 3 harf sıra bilgisini, 270 Form-id bilgisini ifade etmektedir,
- (c): 8 harf sıra bilgisini, 6 Form-id bilgisini ifade etmektedir,
- (ç): 17 harf sıra bilgisini, 1299 Form-id bilgisini ifade etmektedir,
- (d): 23 harf sıra bilgisini, 14586 Form-id bilgisini ifade etmektedir.

MFHD-U (MFHD-Upper Case): MFHD veri kümesinde bulunan büyük harfleri içermektedir. Türk alfabesinde bulunan 'A' dan 'Z' ye büyük harflerden oluşmaktadır. 1'den 29'a kadar olan sayılarla, Türk alfabesinde bulunan 29 büyük harf sırasına göre eşleştirilmiştir. Örneğin; A harfi '1' ile Z harfi '29' ile etiketlenmiştir. MFHD-U veri kümesine ait etiketlenmiş rastgele karakter örnekleri Şekil 3.19'de gösterilmiştir.

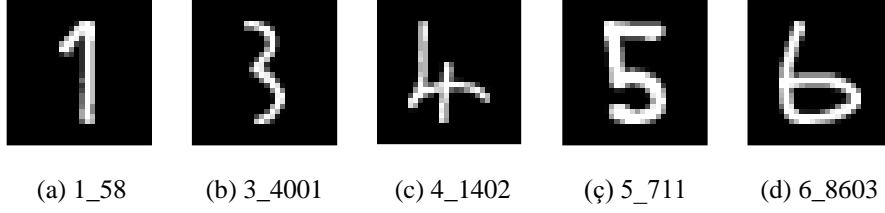


Şekil 3.19 Etiketlenmiş rastgele büyük harf karakter örnekleri.

Şekil 3.19 Açıklaması:

- (a): 2 harf sıra bilgisini, 3 Form-id bilgisini ifade etmektedir,
- (b): 6 harf sıra bilgisini, 1453 Form-id bilgisini ifade etmektedir,
- (c): 9 harf sıra bilgisini, 413 Form-id bilgisini ifade etmektedir,
- (ç): 13 harf sıra bilgisini, 63 Form-id bilgisini ifade etmektedir,
- (d): 17 harf sıra bilgisini, 5198 Form-id bilgisini ifade etmektedir.

MFHD-D (MFHD-Digits): MFHD veri kümesinde bulunan rakamları içermektedir. Her rakama kendi numarası etiket olarak verilmiştir. MFHD-D veri kümesine ait etiketlenmiş karakter örnekleri Şekil 3.20'de gösterilmiştir.



Şekil 3.20 Etiketlenmiş rastgele rakam karakter örnekleri.

Şekil 3.20 Açıklaması:

- (a): 1 rakam sıra bilgisini, 58 Form-id bilgisini ifade etmektedir,
(b): 3 rakam sıra bilgisini, 4001 Form-id bilgisini ifade etmektedir,
(c): 4 rakam sıra bilgisini, 1402 Form-id bilgisini ifade etmektedir,
(ç): 5 rakam sıra bilgisini, 711 Form-id bilgisini ifade etmektedir,
(d): 6 rakam sıra bilgisini, 8603 Form-id bilgisini ifade etmektedir.

MFHD dosya yapısı olarak 4 kısımdan oluşmaktadır. 1. kısımda MFHD-Feature.csv (nitelik bilgilerini içeren dosya, Bakınız Şekil 3.4) dosyası, 2. kısımda MFHD-L (Küçük harf) kümesi, 3. kısımda MFHD-U (Büyük harf) kümesi ve 4. kısımda MFHD-D (Rakam) kümesi bulunmaktadır. MFHD veri kümesi MFHD-L, MFHD-U ve MFHD-D olarak etiketlenmiş ve araştırmacılarla bu şekilde paylaşılmıştır.

4. ARAŞTIRMA BULGULARI ve TARTIŞMA

Bu bölümde MFHD veri kümesi üzerinde farklı yöntemler kullanılarak el yazısı karakter tanıma performansı ölçülmüş ve el yazısı rakam tanıma için kullanılan MNIST ile sonuçların karşılaştırması yapılmıştır.

Deneylerimiz iki kısımdan oluşmaktadır:

- i) MFHD-L, MFHD-U ve MFHD-D veri kümeleri ile karakterlerin sınıflandırılması ve karakterlerden nitelik bilgilerinin çıkarılması,
- ii) MFHD-D rakam tanıma performansının MNIST ile karşılaştırmasının yapılması.

Birinci kısımda yapılan deneylerde literatürdeki çalışmalar dikkate alındığında ML olarak sıkça kullanılan KNN (Fix ve Hodges, 1952) ve DL olarak da LeNet5 (Lecun ve ark., 1998) yöntemi kullanılmıştır. KNN, 28x28 piksellik iki seviyeli (bilevel) görüntüleri giriş olarak almakta ve mesafe metriklerine göre karakter sınıflandırma işlemi yapmaktadır.

MFHD-L, MFHD-U ve MFHD-D ile yapılan deneylerde; Oduntan ve ark. (2018) tarafından yapılan metrik karşılaştırma çalışmasında daha iyi sonuç veren kosinüs mesafe (cosine distance) metriği sınıflandırma çalışmamızda kullanılmıştır. Kosinüs mesafe metrik formülü denklem 4.1 (Fix ve Hodges, 1952)'de gösterilmiştir.

$$k(x, y) = \frac{xy^T}{\|x\| \|y\|} \quad (4.1)$$

LeNet5 ile yapılan bütün sınıflandırma deneylerimizde batch size 128, epochs 30 olarak alınmıştır. LeNet5 mimarisi 32x32 piksellik görüntüleri giriş olarak kabul

etmektedir ancak MFHD veri kümesi MNIST gibi 28x28 pikseldir. Bunun üstesinden gelmek için görüntü çevresi sınırlarla doldurulmuştur (Lecun ve ark., 1998).

İkinci kısımda MFHD-D ve MNIST üzerinde kosinüs mesafe (Cosine distance) metriği ile KNN kullanılarak sıkça kullanılan rakam tanıma yöntemlerinin performansları karşılaştırılmıştır. Aynı şekilde DL yöntemi olarak da LeNet5 (Lecun ve ark., 1998), Effective CNN (Bharadwaj ve ark., 2020) ve Simple CNN (An ve ark., 2020) modelleri kullanılmıştır.

Grover ve Toghi (2019), KNN sınıflandırmada kullanılan k metriğinin optimum değerini sınıflandırma sonuçlarına göre 3, 5 ve 7 olarak önerdiklerinden biz de MFHD-D ve MNIST ile birlikte yapılan deneylerde bu k değerlerini kullandık.

MFHD-L, MFHD-U, MFHD-D ve MNIST veri kümeleri üzerinde KNN, LeNet5 (Lecun ve ark., 1998) ve Effective CNN (Bharadwaj ve ark., 2020) ile yapılan deneyler Google Colab ortamında (25.46 GB RAM, 166.83 GB Disk ve NVIDIA TESLA P100 GPU) gerçekleştirilmiştir. Simple CNN (An ve ark., 2020) mimarisinin eğitimi sırasında epochs sayısının büyüklüğünden dolayı Google Colab üzerinde oluşabilecek kesintilerin önüne geçmek için yerel makinede (Core i5-9300H 2.40GHz, 16 GB DDR4 ve GPU NVIDIA 3 GB GDDR5 GeForce GTX 1050) çalışılmıştır.

4.1. MFHD Üzerinde Yapılan Deneyler

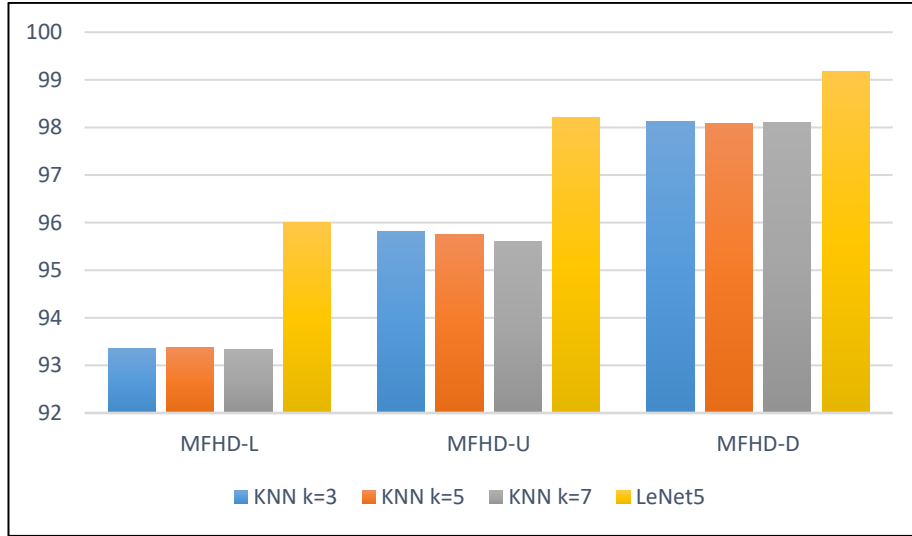
Bu bölümde MFHD veri kümesinde bulunan küçük harf (MFHD-L), büyük harf (MFHD-U) ve rakamlar (MFHD-D) üzerinde KNN ve LeNet5 mimarisi kullanılarak el yazısı karakter sınıflandırma işlemleri yapılmıştır. Öncelikle küçük harf (29 Sınıf), büyük harf (29 Sınıf) ve rakamlar (10 sınıf) üzerinde ayrı ayrı sınıflandırma yapılmıştır. Daha sonra her küçük harf, büyük harf ve rakam için TTT formu ile toplanan 4 farklı nitelik türüne (cinsiyet, eğitim, yaş grubu ve hobi) göre sınıflandırmaları yapılmıştır.

4.1.1. MFHD-L, MFHD-U ve MFHD-D ile sınıflandırma

MFHD-L, MFHD-U ve MFHD-D veri kümeleri kullanılarak KNN ve LeNet5 ile el yazısı karakter sınıflandırma işlemi yapılmıştır. MFHD-L ve MFHD-U veri kümeleri, 29 sınıftan ve 580 000 karakter görüntüsünden, MFHD-D veri kümesi, 10 sınıftan ve 200 000 karakter görüntüsünden oluşmaktadır. MFHD-L, MFHD-U ve MFHD-D veri kümelerinin %80'i eğitim verisi, %20'si test verisi olarak ayrılmıştır. Eğitim kümesi üzerinde uygun parametrelerin belirlenmesi için 5 katlı (5 folds) çapraz doğrulama (Cross Validation) uygulanmıştır. MFHD-L ve MFHD-U veri kümelerinin eğitim verisinde 464 000, test verisinde 116 000 görüntü bulunmaktadır. MFHD-D veri kümesinin eğitim verisinde 160 000, test verisinde 40 000 görüntü bulunmaktadır. MFHD-L, MFHD-U ve MFHD-D için eğitim ve test kümeleri belirlenirken nitelik bilgileri (cinsiyet, eğitim, yaş grubu ve hobi) dikkate alınarak orantılı bir şekilde karakterlerin dağıtımı yapılmıştır. MFHD-L, MFHD-U ve MFHD-D veri kümeleri üzerinde KNN ve LeNet5 yöntemleri yapılan el yazısı karakter sınıflandırma sonuçları tablo olarak Çizelge 4.1'de, grafik olarak da Şekil 4.1'de gösterilmiştir.

Çizelge 4.1 KNN ve LeNet5 yöntemleri ile MFHD (MFHD-L, MFHD-U, MFHD-D) sınıflandırma sonuçları

Veri Kümeleri	Modeller	k Değeri	Test Doğruluğu (%)
MFHD-L	KNN	3	93.34
		5	93.37
		7	93.34
	LeNet5	-	96.00
MFHD-U	KNN	3	95.82
		5	95.74
		7	95.60
	LeNet5	-	98.21
MFHD-D	KNN	3	98.12
		5	98.09
		7	98.11
	LeNet5	-	99.18



Şekil 4.1 KNN ve LeNet5 yöntemleri ile MFHD (MFHD-L, MFHD-U, MFHD-D) sınıflandırma sonuçları

MFHD üzerinde yapılan sınıflandırma sonuçları (Çizelge 4.1 ve Şekil 4.1) incelendiğinde; üç veri kümesi üzerinde en iyi performans LeNet5 için elde edilmiştir. Aynı zamanda, MFHD-D veri kümesi ile yapılan rakam tanıma sınıflandırma performansının MFHD-L ve MFHD-U veri kümelerine nazaran daha iyi olduğu görülmüştür. Bu durum tartışma kısmında ele alınmıştır.

4.1.2. MFHD'nin niteliklere göre sınıflandırılması

El yazısı karakterlerden, cinsiyet ve eğitim durumunu belirlemek için MFHD-L, MFHD-U ve MFHD-D veri kümeleri kullanılmıştır. Aynı şekilde yaş ve hobi nitelik bilgilerini elde etmek için de sırasıyla MFHD-LA, MFHD-UA, MFHD-DA ve MFHD-LH, MFHD-UH, MFHD-DH örneklem kümeleri kullanılmıştır. Karakterler; cinsiyete göre 2 sınıftan (bay, bayan), eğitim durumuna göre 4 sınıftan (ilkokul, ortaokul, lise, yüksekokul) oluşmaktadır. Benzer şekilde yaş grubuna göre 4 sınıftan (5-11 yaş arası, 12-19 yaş arası, 20-30 yaş arası, 31-65 yaş arası) ve hobi durumuna göre 8 sınıftan (kitap, TV, Internet, oyun, spor, müzik, resim, gezi) oluşmaktadır.

Cinsiyet ve eğitim durumuna göre sınıflandırma yapılırken nitelik sayıları eşit olduğundan her karakter için 20 000 görüntü kullanılmıştır. Ancak yaş grubu ve hobi

durumuna göre eşit dağılım olmadığı için en küçük örneklem baz alınarak sınıflandırma yapılmıştır. Yaş grubunda en küçük örneklem 31-65 yaş arası sınıfta 559 katılımcı ile olduğundan sınıflandırma için toplam 2 236 görüntü, hobi durumunda en düşük örneklem Resim sınıfında 766 katılımcı ile olduğundan sınıflandırma için toplam 6 128 görüntü kullanılmıştır. Sınıflandırma işlemlerinde %80 eğitim, %20 test oranı uygulanmıştır.

KNN yöntemi ile nitelik sınıflandırma kullanılacak k değeri için 3, 5 ve 7 değerleri yapılan deneylerle ayrı ayrı incelenmiş ve en iyi sonuç k=5 değeri ile alınmıştır. Bu kısımda gösterilen sonuçlar (el yazısı karakterlerin nitelik sınıflandırması) KNN yönteminde k=5 değeri kullanılarak elde edilmiştir.

MFHD-L veri kümesi üzerinde KNN ve LeNet5 modelleri kullanılarak her bir küçük harf için cinsiyet ve eğitim durumuna göre sınıflandırma sonuçları çizelge 4.2'de gösterilmiştir. Yaş grubuna göre sınıflandırma işlemi MFHD-LA örneklem kümesi üzerinde gerçekleştirilmiş ve sonuçları çizelge 4.3'te verilmiştir. Hobi durumuna göre MFHD-LH örneklem kümesi kullanılmış ve sınıflandırma sonuçları çizelge 4.4'te sunulmuştur.

Çizelge 4.2 MFHD-L veri kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak cinsiyet ve eğitim durumuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama

KR	CİNSİYET (%)				EĞİTİM (%)							
	BAY		BAYAN		İLKOKUL		ORTAOKUL		LİSE		YÜKSEKOKUL	
	KNN	LENET5	KNN	LENET5	KNN	LENET5	KNN	LENET5	KNN	LENET5	KNN	LENET5
a	47,65	52,90	65,05	56,80	51,40	41,40	31,90	40,00	32,60	28,70	27,50	30,50
b	43,70	48,75	65,45	59,50	59,90	52,60	30,20	31,70	30,40	29,80	33,70	44,50
c	46,40	46,70	60,65	62,70	47,80	39,40	35,20	33,70	30,50	46,90	43,30	51,00
ç	50,10	56,65	61,90	55,45	58,60	48,00	32,00	31,40	34,60	34,60	36,90	49,20
d	45,75	52,85	65,55	59,30	59,50	49,80	33,60	33,50	32,80	31,20	40,80	45,50
e	44,25	56,10	61,70	52,70	51,40	42,90	31,40	29,50	31,70	34,40	45,30	44,00
f	45,30	50,85	63,25	59,65	57,30	49,10	34,10	40,40	35,00	22,10	30,50	56,70
g	44,70	55,05	62,95	53,00	55,40	45,60	29,40	31,60	32,30	25,00	32,60	50,20
ğ	44,80	46,70	63,45	60,15	54,40	47,30	29,70	34,60	36,30	27,80	36,10	51,10
h	43,15	60,50	60,90	46,10	50,00	42,70	33,10	26,30	28,50	35,80	32,50	38,70
ı	46,15	45,40	59,40	63,75	42,00	43,70	31,90	44,00	27,20	11,70	27,50	36,60
i	47,15	59,85	58,95	48,80	49,30	31,40	34,60	33,60	27,90	33,40	19,90	38,90
j	51,25	56,00	60,15	51,85	61,00	34,00	30,80	25,40	28,40	42,50	25,60	45,90
k	45,70	52,30	62,45	57,55	55,50	56,00	33,80	29,40	32,50	37,30	45,90	47,00
l	44,60	62,20	58,90	42,10	42,50	38,30	38,90	40,10	29,60	43,20	30,80	37,70
m	46,40	59,25	65,70	53,05	39,00	35,50	33,50	31,90	34,70	30,20	29,20	42,80
n	47,15	57,15	63,20	50,80	41,10	35,30	32,00	27,70	32,00	29,10	29,60	47,10
o	43,45	57,65	65,75	57,15	39,10	29,40	30,60	32,10	27,40	33,20	31,20	42,70
ö	44,85	61,90	67,65	49,85	51,10	44,70	34,50	24,80	36,20	33,50	27,00	51,20
p	43,00	52,25	62,70	51,90	49,10	48,20	31,90	28,80	30,50	30,80	29,90	37,40
r	49,45	58,90	59,70	50,85	55,50	52,30	30,50	28,40	31,20	28,10	25,80	36,30
s	50,20	55,80	62,15	56,05	55,00	47,60	28,00	42,10	31,70	18,20	42,20	50,30
ş	47,40	53,55	65,45	57,00	58,20	53,00	30,10	33,20	34,20	35,00	40,90	46,70
t	45,65	53,60	62,90	53,20	47,30	40,30	36,00	32,70	32,80	36,50	40,40	45,80
u	48,65	52,20	64,60	58,35	45,80	37,70	30,30	38,20	30,40	28,10	23,00	37,80
ü	48,20	63,90	68,45	51,90	47,90	46,40	30,80	29,50	37,20	27,50	23,50	38,70
v	50,10	55,00	60,10	55,80	45,10	32,20	30,00	32,70	30,20	34,20	24,10	27,10
y	50,95	54,80	61,65	49,75	46,70	42,50	29,30	23,60	34,90	31,90	37,80	43,10
z	48,20	47,80	60,35	58,55	58,30	49,80	27,50	27,40	31,30	38,00	34,90	43,50
AVG	46,70	54,71	62,79	54,61	50,87	43,35	31,92	32,36	31,90	31,68	32,70	43,38

Çizelge 4.3 MFHD-LA örneklem kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak yaş grubuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama

KR	YAŞ GRUBU (%)							
	5-11 YAŞ ARASI		12-19 YAŞ ARASI		20-30 YAŞ ARASI		31-65 YAŞ ARASI	
	KNN	LENET5	KNN	LENET5	KNN	LENET5	KNN	LENET5
a	49,55	47,75	25,23	32,43	30,63	26,13	26,13	26,13
b	53,15	45,05	34,23	27,93	26,13	35,14	29,73	38,74
c	56,76	38,74	42,34	44,14	31,53	36,94	43,24	52,25
ç	56,76	49,55	32,43	36,94	36,94	31,53	36,04	44,14
d	61,26	44,14	35,14	51,35	36,04	27,93	39,64	41,44
e	59,46	51,35	25,23	33,33	33,33	31,53	37,84	46,85
f	72,07	48,65	43,24	38,74	24,32	41,44	36,04	50,45
g	59,46	49,55	26,13	32,43	22,52	28,83	33,33	39,64
ğ	54,95	54,95	35,14	28,83	28,83	41,44	36,04	37,84
h	54,05	36,04	21,62	46,85	33,33	36,94	25,23	27,03
ı	38,74	27,93	36,04	45,05	31,53	45,95	24,32	25,23
i	45,05	36,04	36,94	38,74	24,32	35,14	18,92	29,73
j	64,86	49,55	34,23	31,53	28,83	22,52	34,23	46,85
k	52,25	55,86	35,14	38,74	27,93	29,73	40,54	37,84
l	36,94	53,15	37,84	36,94	27,03	27,03	23,42	27,93
m	36,94	41,44	32,43	36,94	27,93	36,04	39,64	28,83
n	35,14	35,14	27,93	43,24	28,83	27,93	17,12	27,03
o	44,14	29,73	29,73	39,64	26,13	22,52	21,62	43,24
ö	54,05	53,15	31,53	28,83	23,42	34,23	35,14	36,04
p	46,85	44,14	26,13	34,23	38,74	45,95	25,23	30,63
r	48,65	43,24	30,63	45,05	28,83	23,42	28,83	36,94
s	56,76	36,94	35,14	40,54	27,03	37,84	50,45	44,14
ş	49,55	45,05	31,53	36,94	28,83	28,83	45,95	54,05
t	57,66	45,95	48,65	45,05	26,13	36,04	34,23	27,93
u	45,95	35,14	34,23	32,43	28,83	30,63	46,85	36,94
ü	45,95	44,14	37,84	23,42	24,32	33,33	45,95	42,34
v	40,54	43,24	34,23	19,82	30,63	29,73	19,82	38,74
y	43,24	50,45	30,63	25,23	32,43	27,93	28,83	24,32
z	49,55	52,25	35,14	27,93	36,04	40,54	25,23	36,94
AVG	50,70	44,42	33,33	35,97	29,36	32,87	32,74	37,25

Çizelge 4.4 MFHD-LH örneklem kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak hobi durumuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama

KR	HOBİ (%)															
	KİTAP		TV		İNTERNET		OYUN		SPOR		MÜZİK		RESİM		GEZİ	
	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS
a	26,14	10,46	26,14	11,76	15,69	13,07	10,46	12,42	8,50	9,15	18,95	13,07	9,15	14,38	22,22	28,76
b	24,18	16,99	15,69	14,38	7,84	15,03	11,11	7,84	9,80	13,07	22,22	15,03	8,50	13,73	19,61	18,95
c	23,53	5,88	16,99	16,99	15,03	22,88	13,07	15,69	9,80	12,42	16,34	17,65	9,15	18,30	43,79	31,37
ç	23,53	19,61	16,34	11,11	10,46	10,46	7,84	17,65	9,80	13,07	16,99	13,73	13,07	16,34	31,37	36,60
d	28,10	22,22	15,69	13,07	10,46	18,95	8,50	15,69	11,11	11,76	18,95	16,34	11,11	10,46	31,37	30,72
e	22,88	18,95	13,73	12,42	11,11	18,30	12,42	15,69	7,19	11,76	22,22	26,14	13,73	14,38	31,37	22,22
f	24,18	9,80	14,38	11,11	12,42	20,26	11,11	10,46	7,84	11,11	22,22	9,15	12,42	15,03	13,73	16,34
g	25,49	18,30	15,03	7,19	10,46	17,65	9,80	12,42	12,42	11,11	18,95	15,69	11,76	10,46	16,99	16,34
ğ	28,10	15,03	14,38	9,80	13,73	16,99	7,84	19,61	12,42	15,69	22,22	16,99	7,84	14,38	20,26	22,88
h	24,18	14,38	18,95	10,46	11,11	12,42	11,76	12,42	9,15	9,80	14,38	19,61	9,15	18,95	15,69	22,22
ı	22,22	11,76	15,69	21,57	15,69	16,99	9,80	7,19	15,69	15,03	13,73	22,88	8,50	3,92	14,38	11,76
i	24,18	18,30	18,30	11,76	10,46	12,42	13,73	11,76	9,80	11,76	20,92	9,15	9,15	19,61	11,76	18,95
j	34,64	16,99	18,95	13,73	13,07	10,46	13,73	18,95	8,50	10,46	20,26	13,73	13,73	26,14	22,88	29,41
k	25,49	21,57	11,76	11,11	16,99	9,15	13,07	12,42	11,76	13,07	18,30	14,38	7,84	18,30	28,10	33,33
l	26,14	7,84	15,69	15,69	10,46	22,22	7,84	6,54	9,15	13,07	16,34	24,84	11,11	9,80	11,76	24,84
m	23,53	10,46	16,99	14,38	8,50	7,84	9,15	10,46	8,50	9,15	17,65	16,99	11,76	16,99	23,53	33,99
n	27,45	16,34	16,34	9,80	10,46	18,30	10,46	20,26	9,15	11,76	18,30	11,76	11,76	14,38	20,92	33,33
o	22,22	16,99	14,38	7,84	11,76	11,76	5,88	11,11	12,42	13,07	17,65	13,07	15,03	16,99	16,99	17,65
ö	30,07	18,30	23,53	9,15	13,73	14,38	11,76	13,07	11,11	9,80	22,88	18,95	16,99	15,69	17,65	26,80
p	24,84	16,99	24,18	10,46	11,11	12,42	11,11	14,38	9,15	13,73	14,38	11,11	9,15	17,65	17,65	30,07
r	28,76	11,76	11,76	7,84	9,80	11,11	11,76	20,92	8,50	9,15	21,57	24,84	8,50	11,76	13,73	19,61
s	24,18	11,76	14,38	13,73	15,03	16,34	7,19	9,15	12,42	13,73	17,65	10,46	16,99	18,30	31,37	28,76
ş	29,41	18,95	17,65	9,80	11,11	15,03	8,50	16,34	11,76	13,73	17,65	16,99	12,42	20,26	32,68	25,49
t	23,53	13,07	18,95	10,46	18,95	11,11	5,23	10,46	10,46	16,34	19,61	22,22	8,50	9,80	21,57	22,88
u	22,88	11,76	9,80	9,15	20,92	18,30	12,42	12,42	9,80	8,50	17,65	16,99	14,38	13,73	12,42	21,57
ü	24,18	18,30	14,38	7,19	11,11	15,03	9,15	12,42	7,84	8,50	20,26	19,61	12,42	19,61	18,30	14,38
v	25,49	13,07	13,07	6,54	12,42	16,99	8,50	11,11	12,42	11,11	13,73	16,99	14,38	22,22	12,42	27,45
y	25,49	19,61	15,69	11,11	12,42	13,07	11,11	7,19	11,76	11,76	19,61	13,73	9,80	14,38	20,26	33,33
z	29,41	18,95	15,69	13,07	7,19	17,65	9,15	12,42	8,50	17,65	16,99	13,73	13,73	15,03	27,45	18,95
AVG	25,67	15,32	16,36	11,47	12,40	15,05	10,12	13,05	10,23	12,08	18,57	16,41	11,45	15,55	21,46	24,79

Büyük harfler üzerinde cinsiyet ve eğitim durumuna göre sınıflandırma yapmak için MFHD-U veri kümesi kullanılmış ve sonuçları çizelge 4.5'te gösterilmiştir. Aynı şekilde yaş grubu nitelik bilgisini elde etmek için MFHD-UA ve hobi durumu için de MFHD-UH örneklem kümeleri kullanılmış ve sonuçları sırasıyla çizelge 4.6 ve çizelge 4.7'de gösterilmiştir.

Çizelge 4.5 MFHD-U veri kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak cinsiyet ve eğitim durumuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama

KR	CİNSİYET (%)				EĞİTİM (%)							
	BAY		BAYAN		İLKOKUL		ORTAOKUL		LİSE		YÜKSEKOKUL	
	KNN	LENET5	KNN	LENET5	KNN	LENET5	KNN	LENET5	KNN	LENET5	KNN	LENET5
A	42,80	45,50	62,75	60,45	56,60	43,30	31,20	34,50	36,20	31,20	33,20	44,30
B	44,15	57,45	65,50	55,00	70,90	62,80	30,20	37,80	29,00	40,20	47,10	44,70
C	46,85	48,25	63,25	64,65	46,30	48,00	30,80	32,40	33,00	32,20	47,30	46,20
Ç	47,65	56,75	63,60	53,95	52,40	44,40	33,40	40,10	33,70	32,20	46,30	53,30
D	49,85	59,85	63,40	51,55	62,70	71,10	31,50	30,30	28,20	25,70	44,30	52,10
E	45,40	58,60	64,70	53,25	61,60	54,70	30,00	35,10	33,40	35,90	46,80	38,30
F	46,80	53,50	63,65	54,10	59,50	51,00	32,40	32,00	34,40	36,50	39,90	47,80
G	47,15	56,10	62,60	55,60	48,00	48,60	33,50	26,70	30,50	26,00	37,90	55,00
Ğ	47,30	49,80	62,80	61,10	48,40	44,50	28,20	34,50	32,60	31,20	47,50	44,70
H	41,30	54,70	65,50	49,30	44,70	45,30	34,90	25,10	31,90	32,30	34,90	45,30
I	45,55	52,90	59,40	54,30	43,30	42,00	33,50	37,90	23,30	28,50	33,70	35,90
İ	46,60	46,00	58,40	58,15	50,40	34,60	31,80	28,10	31,50	36,10	28,70	40,40
J	48,20	50,25	61,35	56,30	60,50	50,40	31,90	28,00	29,20	42,90	37,20	35,90
K	46,80	52,50	62,55	57,45	56,80	34,60	32,00	41,90	31,90	28,80	52,90	53,30
L	49,75	49,65	60,35	61,55	57,60	48,30	30,30	33,40	31,80	31,60	50,50	48,10
M	51,85	63,45	60,55	45,25	53,50	47,60	31,50	33,90	33,20	36,00	44,50	48,30
N	52,85	48,45	59,65	62,35	62,10	57,80	27,80	31,00	38,10	35,60	42,00	49,10
O	48,40	58,20	63,15	52,60	34,30	34,80	30,60	13,60	26,00	40,90	35,50	43,60
Ö	44,50	65,10	66,10	48,65	47,80	49,50	32,80	31,30	34,40	25,20	35,90	59,90
P	44,20	58,20	63,40	49,15	52,70	50,20	34,50	29,90	27,80	31,10	39,40	51,30
R	45,40	53,10	62,15	57,95	60,30	46,40	31,50	39,50	27,10	37,10	52,20	50,50
S	50,10	51,10	62,80	61,30	52,40	51,00	29,70	39,50	33,30	24,20	49,60	46,10
Ş	47,95	51,20	65,15	60,90	55,20	52,80	31,70	30,60	34,20	29,60	51,00	56,20
T	43,75	48,65	61,45	58,85	56,30	40,10	30,80	37,40	31,00	28,20	42,60	57,70
U	48,75	59,00	60,45	49,55	50,40	48,50	31,20	25,10	34,80	41,60	36,60	39,60
Ü	47,75	58,05	64,65	53,80	55,40	47,00	30,60	30,10	35,50	29,10	31,00	45,80
V	50,75	55,70	61,85	57,25	45,70	49,70	30,20	29,40	29,90	22,60	32,80	50,20
Y	45,95	53,90	63,35	56,00	58,40	51,10	27,30	24,70	30,60	30,80	49,50	56,30
Z	44,95	49,75	62,15	58,65	49,10	48,60	26,50	25,10	30,10	23,90	37,20	54,10
AVG	47,01	53,99	62,64	55,83	53,56	48,23	31,16	31,69	31,61	31,97	41,66	48,07

Çizelge 4.6 MFHD-UA örneklem kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak yaş grubuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama

KR	YAŞ GRUBU (%)							
	5-11 YAŞ ARASI		12-19 YAŞ ARASI		20-30 YAŞ ARASI		31-65 YAŞ ARASI	
	KNN	LENET5	KNN	LENET5	KNN	LENET5	KNN	LENET5
A	42,34	44,14	36,04	42,34	37,84	37,84	39,64	47,75
B	78,38	48,65	34,23	45,05	34,23	36,04	43,24	54,95
C	49,55	50,45	40,54	39,64	38,74	41,44	42,34	55,86
Ç	57,66	50,45	38,74	29,73	27,93	26,13	45,95	47,75
D	58,56	38,74	36,04	50,45	25,23	36,04	59,46	62,16
E	70,27	56,76	36,04	36,04	31,53	26,13	45,95	45,95
F	54,95	55,86	36,94	44,14	30,63	37,84	43,24	36,94
G	44,14	41,44	43,24	40,54	37,84	34,23	39,64	45,95
Ğ	43,24	38,74	42,34	33,33	32,43	36,04	54,95	41,44
H	45,05	40,54	33,33	39,64	21,62	26,13	47,75	40,54
I	37,84	36,94	24,32	47,75	38,74	36,04	30,63	38,74
İ	53,15	45,05	28,83	28,83	35,14	34,23	28,83	35,14
J	63,96	66,67	27,03	19,82	37,84	33,33	33,33	51,35
K	51,35	52,25	33,33	37,84	36,94	32,43	58,56	56,76
L	61,26	36,04	29,73	50,45	43,24	36,94	50,45	48,65
M	54,95	44,14	35,14	42,34	35,14	26,13	50,45	63,06
N	53,15	45,05	38,74	38,74	40,54	39,64	50,45	54,05
O	40,54	30,63	32,43	37,84	27,03	24,32	28,83	44,14
Ö	48,65	50,45	30,63	44,14	27,93	22,52	38,74	36,94
P	54,05	50,45	27,93	31,53	26,13	45,05	45,05	45,05
R	57,66	43,24	35,14	45,05	29,73	37,84	45,95	40,54
S	53,15	46,85	44,14	36,04	27,93	38,74	48,65	54,05
Ş	52,25	56,76	36,04	44,14	29,73	32,43	50,45	42,34
T	49,55	43,24	32,43	31,53	36,04	39,64	32,43	50,45
U	52,25	38,74	29,73	39,64	27,93	33,33	53,15	49,55
Ü	46,85	51,35	27,03	17,12	27,93	27,03	35,14	45,05
V	47,75	26,13	40,54	38,74	33,33	36,04	42,34	38,74
Y	54,95	51,35	30,63	44,14	34,23	27,03	51,35	46,85
Z	50,45	45,05	22,52	36,94	30,63	29,73	36,04	45,05
AVG	52,69	45,73	33,92	38,40	32,56	33,46	43,90	47,10

Çizelge 4.7 MFHD-UH örneklem kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak hobi durumuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama

KR	HOBİ (%)															
	KİTAP		TV		İNTERNET		OYUN		SPOR		MÜZİK		RESİM		GEZİ	
	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS
A	30,72	10,46	9,80	12,42	7,19	9,15	8,50	18,95	4,58	17,65	16,34	11,11	9,80	12,42	19,61	25,49
B	22,22	18,30	18,30	9,15	9,15	15,69	14,38	22,22	5,88	3,92	16,99	13,73	7,84	7,84	28,76	30,72
C	21,57	14,38	18,30	9,80	12,42	16,34	10,46	19,61	5,23	7,19	16,34	12,42	14,38	14,38	38,56	48,37
Ç	22,22	13,07	10,46	15,03	15,69	14,38	15,03	16,99	8,50	11,76	16,99	21,57	16,34	15,03	30,07	30,72
D	26,80	11,11	14,38	16,99	7,19	9,15	7,19	15,03	11,76	11,76	21,57	12,42	11,11	11,76	35,95	38,56
E	25,49	13,07	13,07	10,46	13,07	12,42	15,69	10,46	8,50	9,15	18,30	16,99	14,38	15,69	28,10	32,68
F	28,10	17,65	19,61	10,46	9,15	19,61	7,84	17,65	8,50	20,26	18,30	22,88	13,73	13,73	13,07	23,53
G	27,45	19,61	16,99	12,42	8,50	11,76	12,42	13,73	9,15	15,69	17,65	17,65	17,65	13,07	18,95	30,07
Ğ	25,49	15,03	9,80	13,07	13,07	14,38	9,80	8,50	5,88	7,84	17,65	18,95	10,46	12,42	24,18	25,49
H	24,18	12,42	18,95	11,76	9,80	16,34	11,76	16,34	7,84	7,84	15,69	17,65	11,11	11,76	17,65	14,38
I	24,18	15,69	19,61	10,46	11,11	10,46	9,80	13,07	8,50	7,84	15,03	18,30	5,23	8,50	13,73	21,57
İ	22,88	17,65	18,30	13,73	9,15	9,80	15,69	30,07	8,50	5,88	15,03	13,07	11,11	7,19	11,76	22,88
J	27,45	10,46	12,42	6,54	11,11	3,27	11,11	13,07	8,50	14,38	20,92	23,53	7,19	10,46	18,95	35,29
K	24,84	17,65	15,69	12,42	11,11	11,76	11,11	11,11	13,73	7,84	27,45	18,95	11,11	10,46	39,22	41,18
L	27,45	18,95	15,69	9,15	9,80	11,76	11,76	13,73	7,19	10,46	18,95	17,65	7,19	9,80	41,18	41,83
M	22,88	10,46	16,99	11,11	9,80	12,42	11,76	15,69	10,46	9,15	20,92	10,46	11,76	15,69	28,76	35,29
N	26,80	11,76	14,38	11,11	12,42	9,15	9,80	12,42	9,80	7,19	20,26	20,26	10,46	12,42	28,76	28,76
O	30,72	12,42	23,53	11,11	8,50	11,11	7,19	7,84	10,46	12,42	19,61	5,23	11,76	15,69	11,76	17,65
Ö	28,76	16,34	18,30	11,11	15,03	11,11	7,19	12,42	5,23	6,54	18,95	24,18	16,34	14,38	20,92	18,30
P	26,80	8,50	11,11	16,99	11,11	13,07	7,19	13,07	8,50	6,54	18,30	21,57	11,76	7,19	16,99	38,56
R	27,45	18,95	9,80	8,50	11,11	7,84	9,15	16,34	8,50	15,03	22,88	18,30	8,50	13,07	32,03	41,83
S	26,14	13,73	16,99	15,03	13,07	12,42	7,84	16,34	11,11	15,03	15,69	15,03	15,69	17,65	28,76	29,41
Ş	22,22	20,26	11,11	9,15	15,03	13,07	10,46	15,03	10,46	9,80	11,11	20,26	15,69	16,34	27,45	29,41
T	27,45	11,11	14,38	7,84	6,54	9,15	6,54	18,30	13,07	10,46	9,80	12,42	10,46	19,61	28,76	16,99
U	30,72	14,38	16,34	15,03	16,34	18,30	8,50	13,07	7,84	15,03	18,30	10,46	13,73	22,22	13,73	26,14
Ü	32,03	22,88	22,88	15,69	9,15	17,65	6,54	11,11	7,84	9,80	19,61	10,46	16,99	20,92	27,45	30,07
V	26,80	9,15	13,73	14,38	10,46	10,46	7,19	15,69	8,50	9,15	12,42	16,99	9,80	8,50	20,92	30,72
Y	24,84	16,34	15,69	10,46	8,50	14,38	7,84	7,19	9,15	13,73	14,38	21,57	11,11	19,61	33,99	46,41
Z	27,45	13,07	18,30	14,38	15,03	13,73	9,15	13,73	12,42	10,46	17,65	15,03	13,73	16,99	30,07	33,33
AVG	26,28	14,65	15,69	11,92	11,02	12,42	9,96	14,79	8,81	10,68	17,69	16,52	11,95	13,61	25,18	30,54

MFHD-D veri kümesi üzerinde cinsiyet ve eğitim durumuna göre sınıflandırma yapılmış ve sonuçları çizelge 4.8'de gösterilmiştir. Benzer şekilde yaş grubu nitelik bilgisini elde etmek için MFHD-DA ve hobi durumu için de MFHD-DH örneklem kümeleri kullanılarak sınıflandırma işlemi yapılmış ve sonuçları sırasıyla çizelge 4.9 ve çizelge 4.10'da sunulmuştur.

Çizelge 4.8 MFHD-D veri kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak cinsiyet ve eğitim durumuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama

KR	CİNSİYET (%)				EĞİTİM (%)							
	BAY		BAYAN		İLKOKUL		ORTAOKUL		LİSE		YÜKSEKOKUL	
	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS
0	44,70	39,25	61,95	68,50	35,80	41,20	31,10	32,40	29,00	25,10	38,90	37,30
1	48,50	51,10	65,50	61,75	55,40	50,20	33,70	37,30	29,10	30,10	46,00	59,70
2	50,80	49,50	63,65	64,45	51,30	45,50	27,40	30,20	28,90	27,90	44,60	55,20
3	49,10	66,50	64,70	44,60	57,10	52,00	28,20	35,10	34,00	30,00	42,50	52,60
4	48,00	46,50	63,05	62,00	51,30	42,90	25,40	31,50	29,10	28,70	40,70	43,10
5	46,20	60,25	63,95	53,15	50,50	56,40	31,30	30,00	32,10	25,80	46,40	50,40
6	44,70	49,15	62,35	60,70	39,40	33,00	29,10	25,70	30,40	33,70	43,50	42,30
7	45,35	60,75	63,45	48,50	52,30	43,70	26,00	31,40	29,30	28,70	37,00	43,80
8	44,25	54,00	67,15	53,90	47,10	57,80	28,20	34,00	31,30	18,60	40,50	41,70
9	43,45	50,70	67,75	62,25	52,10	40,90	27,40	28,70	30,70	32,30	45,60	42,60
AVG	46,51	52,77	64,35	57,98	49,23	46,36	28,78	31,63	30,39	28,09	42,57	46,87

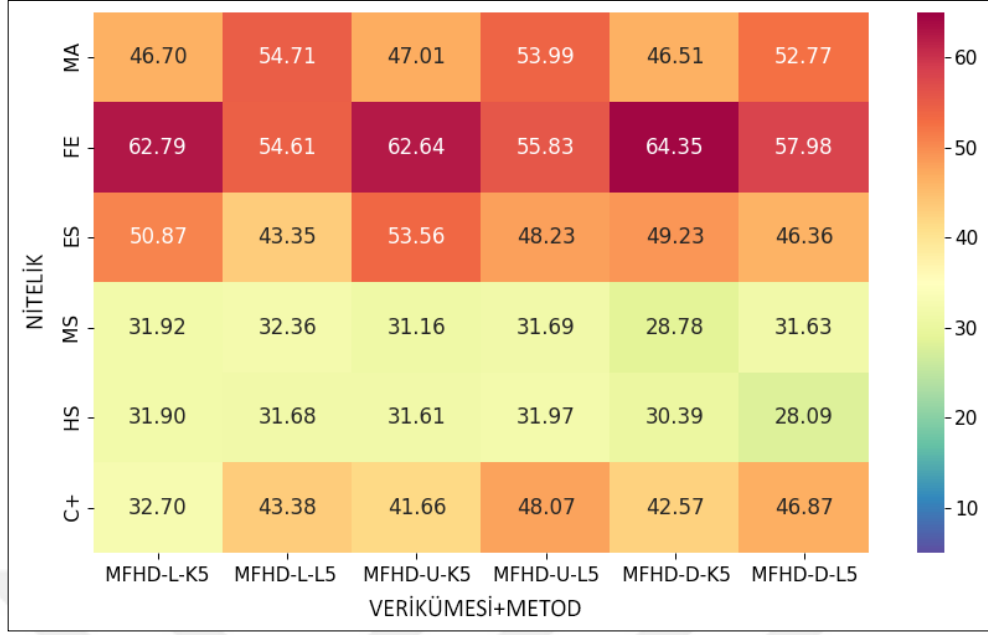
Çizelge 4.9 MFHD-DA örneklem kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak yaş grubuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama

KR	YAŞ GRUBU (%)							
	5-11 YAŞ ARASI		12-19 YAŞ ARASI		20-30 YAŞ ARASI		31-65 YAŞ ARASI	
	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS
0	41,44	34,23	31,53	43,24	36,04	27,93	33,33	36,04
1	58,56	56,76	28,83	36,94	38,74	37,84	48,65	38,74
2	43,24	29,73	24,32	43,24	29,73	22,52	45,95	44,14
3	57,66	60,36	30,63	30,63	31,53	41,44	50,45	53,15
4	41,44	36,04	20,72	38,74	36,94	39,64	41,44	39,64
5	51,35	36,94	27,93	36,94	27,03	27,03	50,45	51,35
6	40,54	43,24	30,63	38,74	30,63	24,32	50,45	54,95
7	48,65	42,34	21,62	33,33	34,23	29,73	39,64	41,44
8	44,14	37,84	30,63	32,43	27,93	28,83	54,05	49,55
9	54,95	38,74	19,82	34,23	36,94	27,03	39,64	37,84
AVG	48,20	41,62	26,67	36,85	32,97	30,63	45,41	44,68

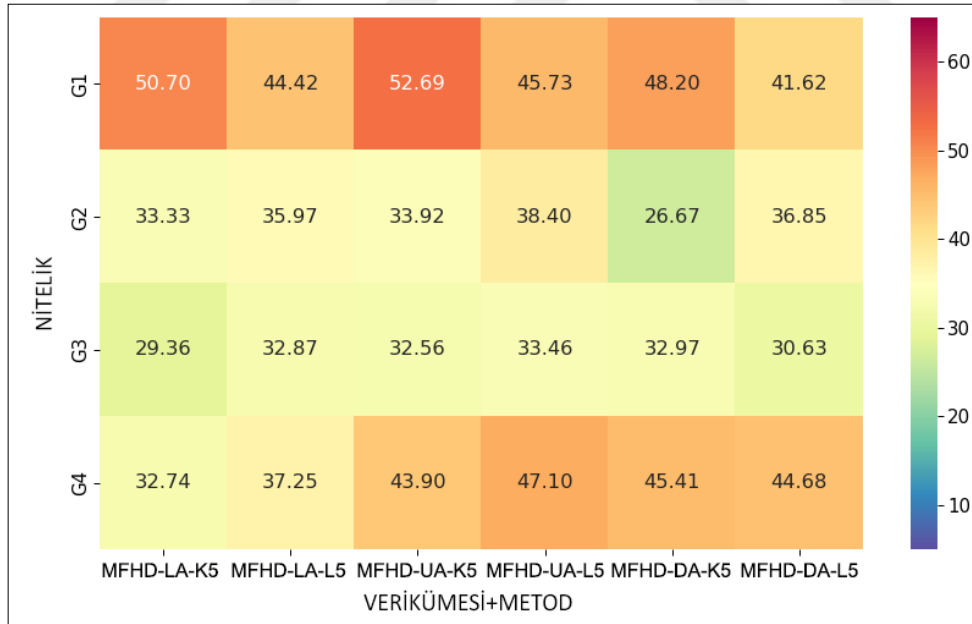
Çizelge 4.10 MFHD-DH örneklem kümesi üzerinde KNN ve LeNet5 yöntemi kullanılarak hobi durumuna göre karakter sınıflandırma sonuçları. Kullanılan kısaltmalar: KR: Karakter, AVG: Ortalama

KR	HOBİ (%)															
	KİTAP		TV		İNTERNET		OYUN		SPOR		MÜZİK		RESİM		GEZİ	
	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS	KNN	LENETS
0	27,45	9,80	21,57	18,95	14,38	7,84	11,11	15,03	7,19	13,73	12,42	7,84	11,11	7,19	9,80	22,22
1	33,33	14,38	13,07	9,80	9,80	16,34	10,46	9,80	8,50	6,54	16,99	24,18	7,19	16,99	28,10	27,45
2	26,14	11,11	7,84	14,38	15,03	17,65	9,15	15,03	13,73	12,42	13,73	13,07	13,07	17,65	16,99	24,18
3	23,53	11,11	11,76	11,76	11,76	14,38	8,50	17,65	7,84	9,15	12,42	18,30	13,07	10,46	37,25	36,60
4	32,03	15,69	15,03	13,73	7,84	16,99	9,15	10,46	12,42	10,46	12,42	11,76	10,46	12,42	22,22	27,45
5	22,88	14,38	11,76	8,50	10,46	13,07	7,19	17,65	7,19	5,88	11,76	12,42	11,76	12,42	24,84	31,37
6	24,84	12,42	14,38	15,03	9,80	16,34	8,50	11,11	7,19	16,99	16,99	8,50	12,42	18,95	18,95	18,30
7	24,18	13,73	18,95	9,80	9,80	15,03	7,84	15,03	11,76	9,15	13,73	15,03	11,11	10,46	20,92	41,83
8	28,76	13,73	8,50	9,15	9,15	20,92	5,88	10,46	9,80	9,15	20,26	15,03	11,76	8,50	23,53	24,18
9	36,60	11,11	11,11	13,07	8,50	13,07	7,19	10,46	11,11	12,42	16,99	20,26	9,80	9,15	26,80	30,07
AVG	27,97	12,75	13,40	12,42	10,65	15,16	8,50	13,27	9,67	10,59	14,77	14,64	11,18	12,42	22,94	28,37

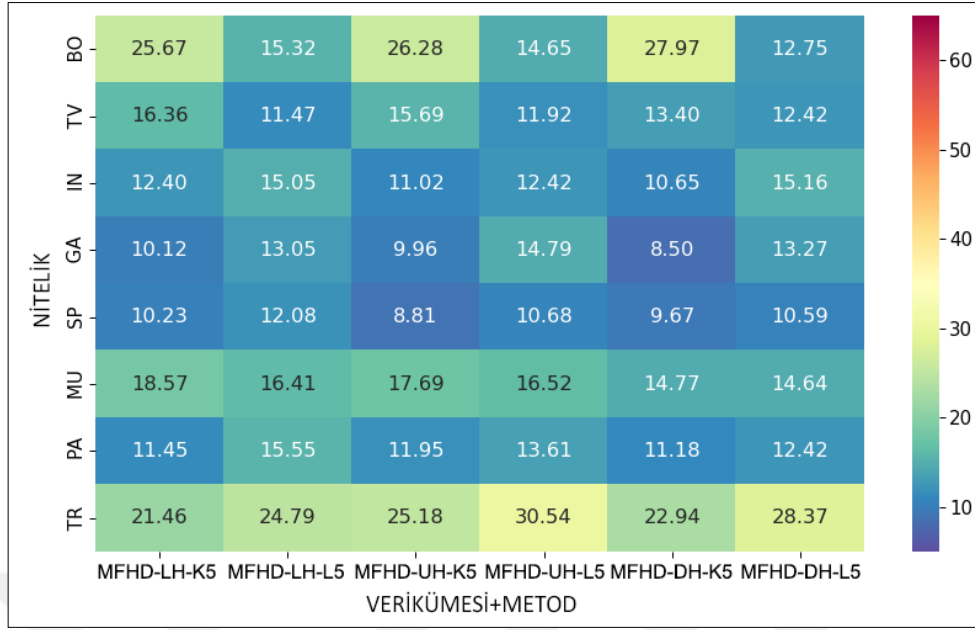
Cinsiyet ve eğitim durumuna göre elde edilen ortalama sınıflandırma doğruluk sonuçlarının (Çizelge 4.2, Çizelge 4.5 ve Çizelge 4.8) sıcaklık haritası gösterimi Şekil 4.2'de, yaş grubuna göre elde edilen ortalama sınıflandırma doğruluk sonuçlarının (Çizelge 4.3, Çizelge 4.6 ve Çizelge 4.9) sıcaklık haritası gösterimi Şekil 4.3'de Hobi durumuna göre elde edilen ortalama sınıflandırma doğruluk sonuçlarının (Çizelge 4.4, Çizelge 4.7 ve Çizelge 4.10) sıcaklık haritası gösterimi Şekil 4.4'te sunulmuştur.



Şekil 4.2 Cinsiyet ve eğitim durumuna göre elde edilen ortalama sınıflandırma doğruluk sonuçlarının (Çizelge 4.2, Çizelge 4.5 ve Çizelge 4.8) sıcaklık haritası gösterimi. Kullanılan kısaltmalar: MA: Bay, FE: Bayan, ES: İlkokul, MS: Ortaokul, HS: Lise, C+: Yüksekokul, K5: KNN (k=5), L5: LeNet5



Şekil 4.3 Yaş grubuna göre elde edilen ortalama sınıflandırma doğruluk sonuçlarının (Çizelge 4.3, Çizelge 4.6 ve Çizelge 4.9) sıcaklık haritası gösterimi. Kullanılan kısaltmalar: G1: 5-11 yaş arası, G2: 12-19 yaş arası, G3: 20-30 yaş arası, G4: 31-65 yaş arası, K5: KNN (k=5), L5: LeNet5



Şekil 4.4 Hobi durumuna göre elde edilen ortalama sınıflandırma doğruluk sonuçlarının (Çizelge 4.4, Çizelge 4.7 ve Çizelge 4.10) sıcaklık haritası gösterimi. Kullanılan kısaltmalar: BO: Kitap, TV: Televizyon, IN: İnternet, GA: Oyun, SP: Spor, MU: Müzik, PA: Resim, TR: Gezi, K5: KNN (k=5), L5: LeNet5

Şekil 4.2’de görüldüğü gibi; cinsiyete göre en yüksek başarı oranı, MFHD-D üzerinde KNN ile yapılan sınıflandırmada bayan niteliğinin tanınmasında çıkmıştır. Aynı şekilde eğitim durumuna göre en yüksek başarı oranı, MFHD-U üzerinde KNN ile yapılan sınıflandırmada ilkokul niteliğinin tanınmasında görülmüştür. Cinsiyete göre en düşük başarı oranı, MFHD-D üzerinde KNN ile yapılan sınıflandırmada bay niteliğinin tanınmasında çıkmış olup benzer şekilde eğitim durumuna göre en düşük başarı oranı, MFHD-D üzerinde LeNet5 ile yapılan sınıflandırmada lise niteliğinin tanınmasında görülmüştür.

Şekil 4.3 incelendiğinde; MFHD-UA üzerinde KNN ile yapılan sınıflandırmada en yüksek başarı oranı G1 (5-11 yaş arası) yaş grubu niteliğinin tanınmasında çıkmıştır. Aynı şekilde en düşük başarı oranı, MFHD-DA üzerinde KNN ile yapılan sınıflandırmada G2 (12-19 yaş arası) yaş grubu niteliğinin tanınmasında görülmüştür.

Şekil 4.4’e baktığımızda; MFHD-UH üzerinde LeNet5 ile yapılan sınıflandırmada en yüksek başarı oranı gezi niteliğinin tanınmasında çıkmıştır. Aynı

şekilde en düşük başarı oranı, MFHD-DH üzerinde KNN ile yapılan sınıflandırmada oyun niteliğinin tanınmasında görülmüştür.

Şekil 4.2, Şekil 4.3 ve Şekil 4.4 birlikte incelendiğinde; niteliklere göre en yüksek başarı oranı, cinsiyete göre yapılan sınıflandırma olurken en düşük başarı oranı da hobi durumuna göre yapılan sınıflandırma olmuştur. Bu durum tartışma kısmında ele alınmıştır.

4.2. MFHD-D ve MNIST ile Yapılan Deneyler

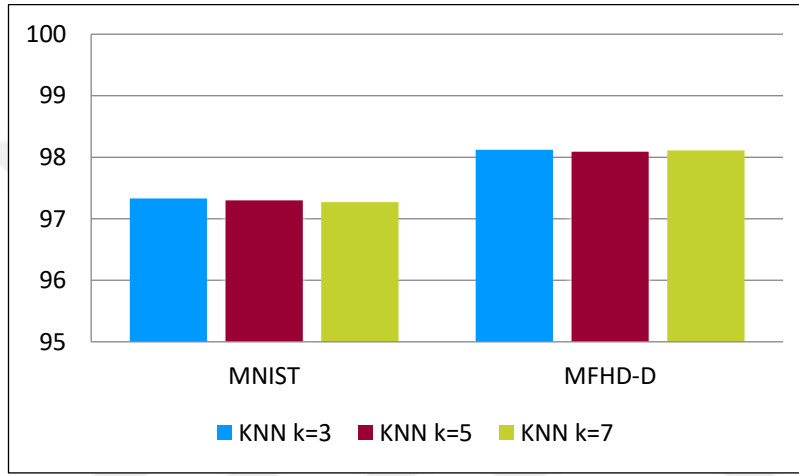
MFHD-D veri kümesinin rakam tanımadaki performansını gözlemlemek için literatürdeki bazı yöntemler (KNN (Fix ve Hodges, 1952), LeNet5 (Lecun ve ark., 1998), Effective CNN (Bharadwaj ve ark., 2020) ve Simple CNN (An ve ark., 2020)) denenmiş ve bu yöntemlerin tutarlılığını belirlemek için sıkça kullanılan MNIST veri kümesi kullanılmıştır. Sınıflandırma için kullandığımız yöntemlerdeki değerler ve parametreler, yazarların çalışmasında belirttiği şekilde uygulanmıştır.

MNIST veri kümesi; 10 sınıf ve 70 000 el yazısı rakamdan (60 000 eğitim, 10 000 test), MFHD-D veri kümesi ise; 10 sınıf ve 200 000 el yazısı rakamdan (160 000 eğitim, 40 000 test) oluşmaktadır. Her iki veri kümesi de yapısal olarak aynı özellikleri (28x28, gri-ton, 8 bit) taşımaktadır. Yaptığımız deneylerde bu eğitim ve test sayıları kullanılmıştır.

MFHD-D ve MNIST veri kümesi üzerinde KNN yöntemi ile yapılan el yazısı rakam tanıma performansları tablo olarak Çizelge 4.11'de, grafik olarak da Şekil 4.5'de gösterilmiştir.

Çizelge 4.11 KNN yöntemi ile MNIST ve MFHD-D sınıflandırma sonuçları

Veri Kümesi	k Değeri	Test Doğruluğu (%)
MNIST	3	97.33
	5	97.30
	7	97.27
MFHD-D	3	98.12
	5	98.09
	7	98.11



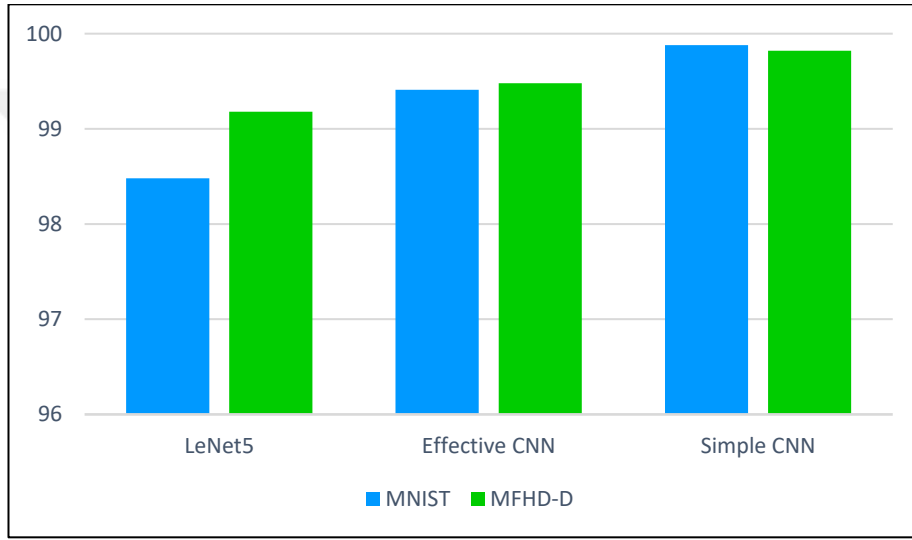
Şekil 4.5 KNN yöntemi ile MNIST ve MFHD-D sınıflandırma sonuçları

Her iki veri kümesi üzerinde yapılan rakam tanıma sınıflandırma sonuçları (Çizelge 4.11 ve Şekil 4.5) incelendiğinde, en iyi performans $k=3$ için elde edilmiştir. KNN yöntemi ile MFHD-D veri kümesi ile yapılan rakam tanıma sınıflandırma performansının MNIST veri kümesine nazaran daha iyi olduğu görülmüştür.

Literatürde MNIST veri kümesi üzerinde denenen bazı yöntemlerin MFHD-D ile tutarlılıklarının incelenmesi yapılmış ve sonuçları tablo olarak Çizelge 4.12'de, grafik olarak da Şekil 4.6'de gösterilmiştir.

Çizelge 4.12 Farklı modeller ile MNIST ve MFHD-D sınıflandırma sonuçları

Veri Kümesi	Modeller	Test Doğruluğu (%)
MNIST	LeNet5 (Lecun ve ark., 1998)	98.48
MFHD-D		99.18
MNIST	Effective CNN (Bharadwaj ve ark., 2020)	99.41
MFHD-D		99.48
MNIST	Simple CNN (An ve ark., 2020)	99.88
MFHD-D		99.82



Şekil 4.6 Farklı modeller ile MNIST ve MFHD-D sınıflandırma sonuçları

MNIST ve MFHD-D veri kümesi üzerinde bazı modeller ile yapılan sınıflandırma sonuçları (Çizelge 4.12 ve Şekil 4.6) incelendiğinde, MFHD-D rakam veri kümesinin eğitim kümesi sayısı ve test kümesi sayısı dikkate alındığında LeNet5 ve Effective CNN modelleri tutarlı sonuç sergilerken Simple CNN'de MNIST veri kümesi daha yüksek doğruluk oranına ulaşmıştır. Bu durum tartışma kısmında ele alınmıştır.

4.3. Tartisima

Bu bölümde; MFHD veri kümesi üzerinde yapılan deneyler analiz edilmiş, kullanılan yöntemlerin performansları karşılaştırılmış ve sık kullanılan veri kümesi MNIST ile sonuçları tartışılmıştır.

Çizelge 4.1 incelendiğinde, MFHD-D üzerinde yapılan deney sonuçlarının, MFHD-L ve MFHD-U veri kümelerine nazaran daha yüksek çıktığı görülmektedir. Bu durumun, MFHD-D kümesinin daha az sınıftan (10 class) oluşması ve birbirine benzeyen karakterler (harflerde olan c, ç, g, ğ, ı, i, o, ö, s, ş, u, ü) içermemesinden kaynaklandığı düşünülmektedir.

MFHD-L, MFHD-U ve MFHD-D veri kümeleri üzerinde yapılan cinsiyet sınıflandırma deneylerinde (Çizelge 4.2, Çizelge 4.5 ve Çizelge 4.8) bayan niteliğinin tanınma oranı, her üç veri kümesinde de ortalamada (AVG) yüksek çıkmıştır. Bu durum; formların katılımcılar tarafından doldurulurken gözlemlenen bayanların baylara nazaran daha özenle yazmasına ve formların elle incelenmesi sırasında bayan yazılarının daha az değişkenlik göstermesine bağlanabilir. Yine bu üç veri kümeleri üzerinde eğitim durumunun analizinde en iyi tanınma olarak ilkökul seviyesi çıkmıştır. Bu sonucun; Erikson'un psikososyal gelişim kuramına göre bireyin gelişim dönemleri dikkate alındığında ilkökul seviyesindeki katılımcıların yazı stillerinin oturmamasına ve formların tek tek elle yapılan incelemelerde karakterlerin kendi içlerinde benzerlik gösterip diğer eğitim seviyelerine göre farklılaşmasına dayandırılabilir. Aynı durum MFHD-LA, MFHD-UA ve MFHD-DA örneklem kümeleri üzerinde yaş grubuna göre yapılan sınıflandırma deneylerinde (Çizelge 4.3, Çizelge 4.6 ve Çizelge 4.9) G1 (5-11 yaş arası) yaş grubunun tanınmasında da gözlemlenmiştir.

Tüm veri kümeleri üzerinde yapılan karakter sınıflandırma ortalamalarına bakıldığında; cinsiyet, eğitim durumu ve yaş grubu tutarlı sonuçlar üretirken hobi durumu değişkenlik göstermiştir. Bu sonucun; katılımcıların el yazısı karakter formunu doldururken cinsiyet, eğitim ve yaş grubu gibi bilgilerini tam olarak

yansıttığı ancak hobi durumunu katılımcının formu doldurduğu andaki ilgisinden hareketle tam olarak yansıtamadığı düşünülmektedir.

Karakter sınıflandırma ortalamalarına bakıldığında; performans başarımları sırasıyla cinsiyet, eğitim durumu, yaş grubu ve hobi durumuna göre olmuştur. Buradan hareketle kişinin cinsiyet bilgisinin el yazısına daha fazla yansıdığı gözlemlenmiştir. Bununla birlikte nitelikli veri kümelerine uygun yöntemlerin geliştirilmesi ile sınıflandırma performansının daha da artacağı ve başarılı sonuçlar üreteceği beklenmektedir.

MNIST VE MFHD-D veri kümeleri üzerinde yapılan deneylerde, karakter tanıma yöntemleri ile sınıflandırma yapılmış ve sonuçları karşılaştırılmıştır. MNIST veri kümesi; 500 katılımcı tarafından yazılmış ve 70 000 karakter görüntüsünden oluşmaktadır. Bundan sebeple sınırlı sayıda katılımcı tarafından yazıldığı için aynı kişilerin elinden çıkmış birden fazla benzer karakter bulunması ve her sınıfta da eşit sayıda karakter bulunmaması el yazısı karakter sınıflandırma modellerinin bazı karakterlere karşı meyilli olmasına yol açabilir. Ancak MFHD-D veri kümesi; 20 000 benzersiz katılımcı tarafından yazılmış, her sınıfta eşit sayıda karakter içeren 200 000 karakter görüntüsünden oluştuğundan daha objektif sınıflandırma yapılmasına imkân sağlayacaktır.

Çizelge 4.11 incelendiğinde MFHD-D'nin tanıma performansının MNIST'ten daha iyi olduğu görülmüştür. Çizelge 4.12'de ise literatürdeki bazı yöntemlerin tutarlı sonuç üretmediği görülmüş, bu durum daha az veri ile öğrenme sürecini gerçekleştirmesine bağlanabilir.

5. SONUÇLAR ve ÖNERİLER

El yazısı tanıma sistemleri geliştirmek ve karşılaştırmak için Türkçe el yazısı karakterlerden oluşan MFHD veri kümesini oluşturmaya yönelik önemli çaba harcanmıştır. MFHD veri kümesi; benzersiz 20 000 katılımcı tarafından nitelikleri ile birlikte, etik ilkelere bağlı kalarak toplanmış, ön işleme adımlarından sonra küçük harf, büyük harf ve rakam türüne göre ayrı ayrı gruplandırılarak etiketlenmiş ve araştırmacıların kullanabileceği bir formata dönüştürülmüştür.

Karakterler, türüne (küçük harf, büyük harf ve rakam) göre sınıflandırılmış ve performansları incelenmiştir. Sık kullanılan MNIST veri kümesi ile MFHD-D rakam veri kümesi farklı metotlar kullanılarak sonuçları karşılaştırılmıştır. Böylelikle hem veri kümelerinin performansları hem de yöntemlerin tutarlılıkları analiz edilmiştir. Niteliklere göre; cinsiyet, eğitim durumu, yaş grubu ve hobi durumuna göre karakterler sınıflandırılmış böylelikle el yazısıyla insan karakterleri ve fizyolojik yapıları arasında bir ilişki olup olmadığı incelenmiştir.

Karakterlerin niteliklere göre sınıflandırma sonuçları incelendiğinde; cinsiyet niteliğinin diğer niteliklere oranla daha iyi ayırt edici olduğu ve başarılı tanıma performansı gösterdiği anlaşılmaktadır. Bu sonuçla kişinin cinsiyet bilgisinin el yazısına daha fazla yansıdığı düşünülmektedir. Nitelikli veri kümelerine uygun yöntemler geliştirilerek daha başarılı sonuçlar üretileceği olası görülmektedir.

Eğitim niteliğini belirlemek için veri kümeleri üzerinde yapılan karakter sınıflandırma deneylerinde, lise seviyesindeki katılımcıların düşük tanınma oranına sahip olması NIST SD 19 (Grother, 1995) veri kümesinde de zorlu bir problem olarak ifade edilmiştir. Bu durumun eğitim bilimleri ve sosyal bilimler için bir araştırma konusu olarak incelenmesinin faydalı olacağını düşünmekteyiz.

El yazısı rakam tanıma çalışmalarında, benzersiz 20 000 katılımcı tarafından yazılan MFHD-D'nin bir karşılaştırma veri kümesi olarak kullanılmasıyla objektif sonuçlar üreteceği düşünülmektedir. Aynı şekilde MFHD veri kümesinin nitelik bilgisi içermesinden dolayı sadece karakter tanıma alanında değil çeşitli bilimsel alanlarda da (Grafoloji, Sosyal ve Beşerî bilimler vb.) kullanılacağı düşünülmektedir.

Türkçe el yazısı karakter kümesi olan MFHD veri kümesinin, Türkçe yazılmış basılı dokümanların dijital ortama aktarılmasına ve tanınmasına önemli bir katkı sağlayacağı düşünülmektedir. Ayrıca ülkemizde yürütülen Eğitim Öğretim faaliyetlerinde ölçme ve değerlendirme yöntemi olan sınavların varlığı göz önüne alındığında MFHD el yazısı veri kümesi kullanılarak derin öğrenme yöntemleri ile klasik sınavların yapılabilmesinin mümkün olabileceği değerlendirilmektedir.

Gelecek çalışmamızda; TTT ile küçük harf ve büyük harf cümle formunda topladığımız pangramlar kullanılarak, derin öğrenme yöntemleri ile bir kişinin el yazısı stili öğrenilebilecektir. Böylelikle el yazısı ile yazılmış dokümanların analizinde nitelik tespitinin daha iyi yapılacağı düşünülmektedir.

KAYNAKLAR

- ABU ALFEILAT, H. A., HASSANAT, A. B., LASASSMEH, O., TARAWNEH, A. S., ALHASANAT, M. B., EYAL SALMAN, H. S., and PRASATH, V. S., 2019. Effects of distance measure choice on k-nearest neighbor classifier performance: a review. *Big data*, 7(4): 221-248.
- AHLAWAT, S., and CHOUDHARY, A., 2020. Hybrid CNN-SVM classifier for handwritten digit recognition. *Procedia Computer Science*, 167: 2554-2560.
- AL-OHALI, Y., CHERIET, M., and SUEN, C., 2003. Databases for recognition of handwritten Arabic cheques. *Pattern Recognition*, 36(1): 111-121.
- AN, S., LEE, M., PARK, S., YANG, H., and SO, J., 2020. An ensemble of simple convolutional neural network models for mnist digit recognition. arXiv preprint arXiv:2008.10400.
- ARICA, N., and YARMAN-VURAL, F. T., 2001. An overview of character recognition focused on off-line handwriting. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 31(2): 216-233.
- ARSLAN, H., and ARSLAN, H., 2021. A new covid-19 detection method from human genome sequences using cpg island features and knn classifier. *Engineering Science and Technology, an International Journal*, 24(4): 839-847.
- ARVANITOPOULOS, N., CHEVASSUS, G., MAGGETTI, D., and SÜSSTRUNK, S., 2017. A handwritten french dataset for word spotting: CFRAMUZ. In *Proceedings of the 4th International Workshop on Historical Document Imaging and Processing* (pp. 25-30).
- BALDOMINOS, A., SAEZ, Y., and ISASI, P., 2019. A survey of handwritten character recognition with mnist and emnist. *Applied Sciences*, 9(15): 3169.
- BARTOS, G. E., HAJNAL, E., and HOŞCAN, Y., 2018. Comparison of Feature Extraction Techniques for Handwriting Recognition. In *2018 IEEE 12th International Symposium on Applied Computational Intelligence and Informatics (SACI)* (pp. 000405-000410), IEEE.
- BARTOS, G. E., HOŞCAN, Y., KAUER, A., and HAJNAL, É. N., 2020. A Multilingual Handwritten Character Dataset: THE Dataset. *Acta Polytechnica Hungarica*, 17(9).
- BAY, H., TUYTELAARS, T., and GOOL, L. V., 2006. SURF: Speeded up robust features. In *European conference on computer vision* (pp. 404-417), Springer, Berlin, Heidelberg.
- BEOHAR, D., and RASOOL, A., 2021. Handwritten digit recognition of MNIST dataset using deep learning state-of-the-art artificial neural network (ANN) and Convolutional Neural Network (CNN). In *2021 International Conference on Emerging Smart Computing and Informatics (ESCI)* (pp. 542-548), IEEE.
- BHARADWAJ, Y. S., RAJARAM, P., SRIRAM, V. P., SUDHAKAR, S., and PRAKASH, K. B., 2020. Effective handwritten digit recognition using deep convolution neural network. *International Journal of Advanced Trends in Computer Science and Engineering*, 9(2): 1335-1339.

- BHATIA, E. N., 2014. Optical character recognition techniques: a review. *International journal of advanced research in computer science and software engineering*, 4(5): 1219-1223.
- BİRİNCİOĞLU, İ., Ö. KURTAŞ, İ. ÇAKIR, N. YURTSEVER, ve H. TEKE, 2010. El yazısı incelemelerinde fulaj kavramı, *Adli Bilimler Dergisi* 9(2): 32–37.
- BOUKERROUI, D., NOBLE, J. A., and BRADY, M., 2004. On the choice of band-pass quadrature filters. *Journal of Mathematical Imaging and Vision*, 21(1): 53-80.
- CAMASTRA, F., SPINETTI, M., and VINCIARELLI, A., 2006. Offline Cursive Character Challenge: a New Benchmark for Machine Learning and Pattern Recognition Algorithms. In *18th International Conference on Pattern Recognition (ICPR'06) (Vol. 2, pp. 913-916)*, IEEE.
- CHA, S. H., 2007. Comprehensive survey on distance/similarity measures between probability density functions. *City*, 1(2): 1.
- CMOLIK, L., PAVLOVEC, V., WU, H. Y., and NOLLENBURG, M., 2020. Mixed labeling: Integrating internal and external labels. *IEEE Transactions on Visualization and Computer Graphics*.
- COHEN, G., AFSHAR, S., TAPSON, J., and VAN SCHAIK, A., 2017. EMNIST: Extending MNIST to handwritten Letters. In *2017 international joint conference on neural networks (IJCNN)*, (pp. 2921-2926), IEEE.
- COVER, T., and HART, P., 1967. Nearest neighbor pattern classification. *IEEE transactions on information theory*, 13(1): 21-27.
- ÇAPAR, A., TAŞDEMİR, K., KILIÇ, Ö., and GÖKMEN, M., 2003. A Turkish handprint character recognition system. In *International Symposium on Computer and Information Sciences* (pp. 447-456), Springer, Berlin, Heidelberg.
- DE CAMPOS, T. E., BABU, B. R., and VARMA, M., 2009. Character recognition in natural images. *VISAPP* (2): 7(2).
- DEMİR, Ü., ve UĞURLU, B., 2021. El Yazısından Kişilik Analizinde El Yazısı Tanılamaya Yönelik Bir Karar Destek Modeli Önerisi. *Uluslararası Batı Karadeniz Mühendislik ve Fen Bilimleri Dergisi*, 3(2).
- DEMİRKAYA, K. G., and ÇAVUŞOĞLU, Ü., 2021. Handwritten Digit Recognition With Machine Learning Algorithms. *Academic Platform Journal of Engineering and Smart Systems*, 10(1): 9-18.
- EL-SAWY, A., LOEY, M., and EL-BAKRY, H., 2017. Arabic handwritten characters recognition using convolutional neural network. *WSEAS Transactions on Computer Research*, 5(1): 11-19.
- EL-SHERIF, E. A., and ABDELAZEEM, S., 2007. A Two-Stage System for Arabic Handwritten Digit Recognition Tested on a New Large Database. In *Artificial intelligence and pattern recognition* (pp. 237-242).
- ERIKSON, E. H., 1993. *Childhood and society*. WW Norton and Company.
- ESPANA, S., CASTRO, M. J., and HIDALGO, J. L., 2004. The SPARTACUS-Database: A Spanish sentence database for offline handwriting recognition. In *Proceedings of the 4th International Conference on Language Resources and Evaluation (LREC 2004)* (pp. 227-230).
- FIX, E., and HODGES Jr, J. L., 1952. Discriminatory analysis-nonparametric discrimination: Small sample performance. *California Univ Berkeley*.

- GOPE, B., PANDE, S., KARALE, N., DHARMALE, S., and UMEKAR, P., 2021. Handwritten digits identification using MNIST database via machine learning models. In IOP Conference Series: Materials Science and Engineering (Vol. 1022: No. 1, p. 012108), IOP Publishing.
- GRIMSDALE, R. L., SUMNER, F. H., TUNIS, C. J., and KILBURN, T., 1959. A system for the automatic recognition of patterns. *Proceedings of the IEE-Part B: Radio and Electronic Engineering*, 106(26): 210-221.
- GROTHER, P. J., 1995. NIST special database 19. Handprinted forms and characters database. National Institute of Standards and Technology, 10.
- GROTHER, P. J., 2017. NIST special database 19. NIST handprinted forms and characters database.
- GROVER, D., and TOGHI, B., 2019. MNIST dataset classification utilizing k-NN classifier with modified sliding-window metric. In *Science and Information Conference* (pp. 583-591), Springer, Cham.
- HAMBAL, A. M., PEI, Z., and ISHABAILU, F. L., 2017. Image noise reduction and filtering techniques. *International Journal of Science and Research (IJSR)*, 6(3): 2033-2038.
- HANSEN, T., and GEGENFURTNER, K. R., 2017. Color contributes to object-contour perception in natural scenes. *Journal of Vision*, 17(3): 14-14.
- HAO, Y., and CHEN, J., 2020. Deep Neural Network for Handwritten Digital Recognition Based on Attention Mechanism.
- HARRIS, C., and STEPHENS, M., 1988. A combined corner and edge detector. In *Alvey vision conference* (Vol. 15: No. 50, pp. 10-5244).
- HEARST, M. A., DUMAIS, S. T., OSUNA, E., PLATT, J., and SCHOLKOPF, B., 1998. Support vector machines. *IEEE Intelligent Systems and their applications*, 13(4): 18-28.
- HUBER, R. A., and HEADRICK, A. M., 1999. *Handwriting identification: facts and fundamentals*. CRC press.
- HULL, J. J., 1994. A database for handwritten text recognition research. *IEEE Transactions on pattern analysis and machine intelligence*, 16(5): 550-554.
- INTERNET, Kişisel Verileri Koruma Kanunu (KVKK), <https://www.kvkk.gov.tr/> (Erişim: 16.12.2019).
- JUSTUSSON, B. I., 1981. Median filtering: Statistical properties. *Two-Dimensional Digital Signal Processing II*, 161-196.
- KAUR, D., and KAUR, Y., 2014. Various image segmentation techniques: a review. *International Journal of Computer Science and Mobile Computing*, 3(5): 809-814.
- KHAN, M. W., 2014. A survey: Image segmentation techniques. *International Journal of Future Computer and Communication*, 3(2): 89.
- KIM, D. H., HWANG, Y. S., PARK, S. T., KIM, E. J., PAEK, S. H., and BANG, S. Y., 1996. Handwritten Korean character image database PE92. *IEICE transactions on information and systems*, 79(7): 943-950.
- KUSETOGULLARI, H., YAVARIABDI, A., CHEDDAD, A., GRAHN, H., and HALL, J., 2020. ARDIS: a Swedish historical handwritten digit dataset. *Neural Computing and Applications*, 32(21): 16505-16518.
- LECUN, Y., BOTTOU, L., BENGIO, Y., and HAFFNER, P., 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11): 2278-2324.

- LECUN, Y., 1998. The MNIST database of handwritten digits. [http://yann. lecun. com/exdb/mnist/](http://yann.lecun.com/exdb/mnist/). (Eriřim: 17.06.2020).
- LI, S., LI, H., LI, M., SHYR, Y., XIE, L., and LI, Y., 2009. Improved prediction of lysine acetylation by support vector machines. *Protein and peptide letters*, 16(8): 977-983.
- LIU, C. L., YIN, F., WANG, D. H., and WANG, Q. F., 2011. CASIA online and offline Chinese handwriting databases. In *2011 International Conference on Document Analysis and Recognition* (pp. 37-41), IEEE.
- LIU, C. L., YIN, F., WANG, D. H., and WANG, Q. F., 2013. Online and offline handwritten Chinese character recognition: benchmarking on new databases. *Pattern Recognition*, 46(1): 155-162.
- LOWE, D. G., 2004. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2): 91-110.
- MAHMOUD, S. A., AHMAD, I., AL-KHATIB, W. G., ALSHAYEB, M., PARVEZ, M. T., MÄRGNER, V., and FINK, G. A., 2014. KHATT: An open Arabic offline handwritten text database. *Pattern Recognition*, 47(3): 1096-1112.
- MAILLO, J., TRIGUERO, I., and HERRERA, F., 2015. A mapreduce-based k-nearest neighbor approach for big data classification. In *2015 IEEE Trustcom/BigDataSE/ISPA* (Vol. 2: pp. 167-172), IEEE.
- MANJUSHA, K., KUMAR, M. A., and SOMAN, K. P., 2019. On developing handwritten character image database for Malayalam language script. *Engineering Science and Technology, an International Journal*, 22(2): 637-645.
- MARTI, U. V., and BUNKE, H., 2002. The IAM-database: an English sentence database for offline handwriting recognition. *International Journal on Document Analysis and Recognition*, 5(1): 39-46.
- MINAEE, S., BOYKOV, Y. Y., PORIKLI, F., PLAZA, A. J., KEHTARNAVAZ, N., and TERZOPOULOS, D., 2021. Image segmentation using deep learning: A survey. *IEEE transactions on pattern analysis and machine intelligence*.
- MOHAMMAD, M. J., REZA, E., and FATEMEH, S., 2011. Persian handwritten digits recognition: A divide and conquer approach based on mixture of MLP experts. *International Journal of Physical Sciences*, 6(30): 7007-7015.
- NURSEITOV, D., BOSTANBEKOV, K., KURMANKHOJAYEV, D., ALIMOVA, A., ABDALLAH, A., and TOLEGENOV, R., 2021. Handwritten Kazakh and Russian (HKR) database for text recognition. *Multimedia Tools and Applications*, 80(21): 33075-33097.
- ODUNTAN, O. E., ADEYANJU, I. A., FALOHUN, A. S., and OBE, O. O., 2018. A comparative analysis of euclidean distance and cosine similarity measure for automated essay-type grading. *Journal of Engineering and Applied Sciences*, 13(11): 4198-4204.
- OLIVEIRA, L. S., and SABOURIN, R., 2004. Support vector machines for handwritten numerical string recognition. In *Ninth international workshop on frontiers in handwriting recognition* (pp. 39-44), IEEE.
- PARE, S., KUMAR, A., SINGH, G. K., and BAJAJ, V., 2020. Image segmentation using multilevel thresholding: a research review. *Iranian Journal of Science and Technology, Transactions of Electrical Engineering*, 44(1): 1-29.

- PASHINE, S., DIXIT, R., and KUSHWAH, R., 2021. Handwritten digit recognition using machine and deep learning algorithms. arXiv preprint arXiv:2106.12614.
- PENG, B., ZHANG, L., and ZHANG, D., 2013. A survey of graph theoretical approaches to image segmentation. *Pattern recognition*, 46(3): 1020-1038.
- PLAMONDON, R., and SRIHARI, S. N., 2000. Online and off-line handwriting recognition: a comprehensive survey. *IEEE Transactions on pattern analysis and machine intelligence*, 22(1): 63-84.
- PUROHIT, A., and CHAUHAN, S. S., 2016. A literature survey on handwritten character recognition. *IJCSIT) International Journal of Computer Science and Information Technologies*, 7(1): 1-5.
- PRATHAP, K. S. V., JILANI, S. A. K., and REDDY, P. R., 2016. A critical review on Image Mosaicing. In *2016 International Conference on Computer Communication and Informatics (ICCCI)* (pp. 1-8), IEEE.
- ROBERTS, L. G., 1963. *Machine perception of three-dimensional soups*. Massachusetts Institute of Technology, 2017.
- ROBERTS, P., 2002. *Love Letters: The Romantic Secrets Hidden in Our Handwriting*. Career Press.
- RUBLEE, E., RABAUD, V., KONOLIGE, K., and BRADSKI, G., 2011. ORB: An efficient alternative to SIFT or SURF. In *2011 International conference on computer vision* (pp. 2564-2571), IEEE.
- SAITO, T., 1985. On the data base ETK9B of handprinted characters in JIS Chinese characters and its analysis. *IEICE trans*, 68(4): 757-772.
- SAXENA, C., and KOURAV, D., 2014. Noises and image denoising techniques: a brief survey. *International journal of Emerging Technology and advanced Engineering*, 4(3): 878-885.
- SEZGIN, M., and SANKUR, B., 2004. Survey over image thresholding techniques and quantitative performance evaluation. *Journal of Electronic imaging*, 13(1): 146-165.
- SMITH, S. M., and BRADY, J. M., 1997. SUSAN—a new approach to low level image processing. *International journal of computer vision*, 23(1): 45-78.
- SOBEL, I., and FELDMAN, G., 1968. A 3x3 isotropic gradient operator for image processing. a talk at the Stanford Artificial Project in, 271-272.
- SRL, T., 1994. The Semeion database of handwritten digits. <http://archive.ics.uci.edu/ml/datasets/Semeion?Handwritten?Digit>. (Erişim: 17.08.2021).
- ŞEKERCİ, M., ve KANDEMİR, R., 2006. Sözlük kullanarak Türkçe el yazısı tanıma. *Elektrik–Elektronik Bilgisayar–Mühendisliği Sempozyum (ELECO 2006)*, Aralık.
- ŞEKERCİ, M., 2007. Birleşik ve eğik Türkçe el yazısı tanıma sistemi (Master's thesis, Trakya Üniversitesi, Fen Bilimleri Enstitüsü).
- THOMA, M., 2017. The hasyv2 dataset. arXiv preprint arXiv:1701.08380.
- TOSELLI, A. H., JUAN, A., GONZÁLEZ, J., SALVADOR, I., VIDAL, E., CASACUBERTA, F., and NEY, H., 2004. Integrated handwriting recognition and interpretation using finite-state models. *International Journal of Pattern Recognition and Artificial Intelligence*, 18(04): 519-539.
- TRAJKOVIĆ, M., and HEDLEY, M., 1998. FAST corner detection, *Image Vis. Comput.* 16 (2): 75–87.

- UĞURLU, B., KAÇAN, K., ve TÜRKYILMAZ, İ., 2010. Bilgi Güvenliğinde El Yazısı. Akademik Bilişim, 10: 409-14.
- VAGHELA, D., and NAINA, P., 2014. A review of image mosaicing techniques. arXiv preprint arXiv:1405.2539.
- VAPNIK, V. N., 1995. Constructing learning algorithms. In The nature of statistical learning theory (pp. 119-166), Springer, New York, NY.
- VURAL, E., ERDOGAN, H., OFLAZER, K., and YANIKOGLU, B., 2004. An online handwriting recognition system for Turkish. In Proceedings of the IEEE 12th Signal Processing and Communications Applications Conference, 2004. (pp. 607-610), IEEE.
- WANG, F., WANG, Q., NIE, F., YU, W., and WANG, R., 2018. Efficient tree classifiers for large scale datasets. Neurocomputing, 284: 70-79.
- WU, H., 2018. CNN-Based Recognition of Handwritten Digits in MNIST Database. Research School of Computer Science. –The Australia National University, Canberra.
- YANG, D., PENG, B., AL-HUDA, Z., MALIK, A., and ZHAI, D., 2022. An overview of edge and object contour detection. Neurocomputing.
- YU, N., JIAO, P., and ZHENG, Y., 2015. Handwritten digits recognition base on improved LeNet5. In The 27th Chinese Control and Decision Conference (2015 CCDC) (pp. 4871-4875), IEEE.
- ZHANG, S., LI, X., ZONG, M., ZHU, X., and WANG, R., 2017. Efficient kNN classification with different numbers of nearest neighbors. IEEE transactions on neural networks and learning systems, 29(5): 1774-1785.