

Zaman Serisi Verilerinde Anomali Tespiti için 5 Modelin Eğitilmesi ve Karşılaştırılması

Furkan YILMAZ & Eyup CINAR

Öz: Anomali tespiti, verilerden aykırı bileşenleri tespit etme ve çıkarma için yıllardır kullanılmaktadır. Anomali tespitinde birçok teknik kullanılmaktadır. Önemi her gün artan tekniklerden biri de bu alanda büyük bir rol oynayan makine öğrenmesidir. Bu çalışmada, Airbus Sas tarafından sağlanan ağır donanımlı bir helikopterin farklı bölgelerine yerleştirilen ivme ölçerlerden elde edilen zaman serisi verisi kullanılarak titreşim seviyelerinin analizi yapılmıştır. Veri, 1 boyutlu 1 dakikalık dizilere dönüştürülmüş ve bu formatta hazırlanmıştır. Beş farklı model (Isolation Forest, One-Class SVM, K-means, GMM ve LSTM) kullanılarak eğitim yapılmış ve veri setine üç farklı imputasyon yöntemi (Ortalama, K-NN ve FB-Fill) uygulanmıştır. Elde edilen sonuçlara göre, 1C-SVM ve LSTM modelleri model eğitimi sonucunda en yüksek F1 skoru olan 0.89'u elde etmiştir. Bu bulgular, helikopter titreşim seviyelerinin etkili bir şekilde modelleyebilen ve anomali tespiti konusunda başarılı sonuçlar veren bu iki modelin öne çıktığını göstermektedir. Çalışma, havacılık endüstrisinde güvenlik ve bakım optimizasyonu açısından önemli bir adım olan titreşim analizi konusunda değerli bilgiler sunmaktadır.

Anahtar kelimeler: anomali tespiti, zaman serisi verileri, makine öğrenmesi.

1. GİRİŞ

Anomali tespiti (ayrıca **aykırı değer** ve bazı yerlerde **yenilik tespiti** olarak adlandırılır) yıllardır üzerinde çalışılan bir konudur. Farklı uygulamalar üzerinde anomali tespiti için çeşitli şahsına münhasır metodlar geliştirilmiş ve kullanılmıştır, [1]. Anomali Tespiti “alışılmışın dışında davranan veri modellerini bulma sorunu” olarak tanımlanır, [2],[3].

Anomali tespitinin uygulama alanları şunlardır; izinsiz giriş tespiti, dolandırıcılık tespiti, arıza tespiti, sistem sağlığı izlenmesi, sensör ağlarındaki olay tespiti, ekosistem bozukluklarının tespiti ve görüntülerde kusur tespiti[5]. Bu çalışmada da sensör verileri ve zaman serisi şeklinde olan veriler üzerinde anomali tespiti yapan makaleler incelenmektedir. Zaman serisi verileri, düzenli zaman aralıkları ardışık olarak ölçülen verilerdir.

Bu araştırmanın kullandığı veri seti Airbus SAS tarafından toplanmış ve yayınlanmıştır. Ağır donanımlı hava araçlarının uçuş testlerindeki ana problem toplanan sinyal verilerinin teyidinin, verilerin çokluğu yüzünden zor olmasıdır. Elle teyit çok fazla zaman ve insan gücüne ihtiyaç duyduğu için bu işlemin otonom hale gelmesi kritiktir. Bu durumda, helikopterin tüm çalışma şartlarındaki titreşim seviyelerini ölçmek için, helikoptere farklı yönlerde (enlemesine, boylamasına, yanlamasına)

ve farklı konumlarda ivmeölçerler konulmuştur. Araştırmada kullanılan veri seti farklı uçuşlardan toplanan 1024 Hz frekansında 1 dakikalık dizilere bölünmüş 1B zaman serilerinden oluşmaktadır [4].

Bu araştırmanın amacı bu veri setini kullanarak anomali tespiti için dört adet makine öğrenmesi ve bir adet derin öğrenme modeli eğitmektir. Makine öğrenmesi modelleri şunlardır; **Isolation Forest, K-Means Cluster, Gaussian Mixture Model, One-Class SVM**. Derin öğrenme modeli ise **LSTM**'dir. Eğitilen modellerin sonuçları karşılaştırılıp anomali tespiti için en iyisi hangisi karar verilecektir.

2. SEMBOLLER, BİRİMLER, KISALTMALAR VE TANIMLAR

İsim	Açıklama
Fuzzy C-Mean Clustering (FCM)	Kümeleme problemlerinde kullanılan bir yöntemdir.
Least Squares Support Vector Machine (LS-SVM)	Regresyon veya sınıflandırma problemleri için kullanılan bir öğrenme modelidir.
Recurrent Neural Networks (RNN)	Zaman serileri ve sıralı verilerde etkili olan bir tür yapay sinir ağıdır.
Long Short-	RNN'lerdeki problemleri

Term Memomry (LSTM)	çözmek için tasarlanmış bir özel hücre tipidir ve uzun vadeli bağımlılıkları daha etkili bir şekilde öğrenir.
Gated Recurrent Unit (GRU)	LTSM'nin daha basitleştirilmiş bir versiyonudur ve benzer performansı sağlarken daha az hesaplama maliyetine sahiptir.
Support Vector Data Description (SVDD)	Sağlıklı örnek verilerin işlenerek tespit eşiği oluşturan sınıflandırıcıdır.
Gauss Karışım Modeli (GMM)	Bir sistemin normal davranışını modellemek ve bu normal davranıştan sapmaları belirlemek için kullanılabilen istatistiksel bir yöntemdir.
Holt-Winters	Zaman serileri analizi için kullanılan bir ön görüleme yöntemidir ve zaman içindeki trend, mevsimsellik ve seviye değişimlerini modellemek amacıyla üç temel bileşeni içerir.
Auto Regressive Integrated Moving Average (ARIMA)	Zaman serileri analizi ve ön görüleme için kullanılan istatistiksel bir modeldir.
Gradient Boosted Regression Trees (GBRT)	Ensemble öğrenme yöntemi olan Gradient Boosting'in bir türüdür ve regresyon problemlerinde güçlü tahminler yapmak için karar ağaçlarını birleştirir.
Long Short-Term Memory-Variational Autoencoder (LSTM-VAE)	Zaman serileri üzerinde değişkenlik analizi için kullanılan bir yapay sinir ağı modelidir. LSTM uzun vadeli bağımlılıkları yakalarken, VAE de değişkenlik analizi ve veri üretimi için kullanılır, böylece birleştirilmiş model zaman serileri verilerini etkili bir şekilde öğrenir ve analiz eder.

Z-Score	Bir veri noktasının ortalama değerden ne kadar uzakta olduğunu, ölçümün standart sapma birimleri cinsinden ifade eden istatistiksel bir ölçüdür.
K-NN Imputation	Eksik değerlere sahip bir gözlem biriminin eksik olan değerini, benzer diğer gözlem birimlerinin bilinen değerleri kullanılarak tahmin etmeye dayanan bir imputasyon yöntemidir.
Forward-Backward Fill	Birbirini takip eden zaman serisi gözlemlerindeki eksik değerleri sırasıyla önceki ve sonraki değerlerle dolduran imputasyon yöntemleridir.
Principal Component Analysis (PCA)	Çok boyutlu bir veri setindeki değişkenler arasındaki ilişkileri anlamak ve veri setini daha az sayıda, bağımsız değişkenle ifade etmek için kullanılan bir boyut indirgeme tekniğidir.

1

3. LİTERATÜR İNCELEMESİ

2

2.1 Makale 1: “Anomaly Detection in Time Series Data using a Fuzzy C-Means Clustering”[5]

3

Bu makalede araştırmacılar **Fuzzy C-Means Clustering** (bulanık kümeleme) kullanarak zaman serisi verisinde anomali tespiti üzerinde çalışma yapmışlardır. Bu algoritma belirli bir veri kümesini belirli özelliklere göre gruplamak için kullanılır.

10

Araştırmacılar ilk olarak zaman serilerindeki anomalileri genlik ve şekil anomalisi olarak ikiye ayırmaktadır. İki tür anomalinin tespitini sağlayacak bir birleşik çerçeve kullanmışlar. Bulanık kümeleme yöntemi kullanılarak zaman serilerindeki yapı ortaya çıkarılmış ve bir yeniden yapılandırma kriteri kullanılarak her bir alt diziye anomali skoru tayin edilmiştir. Genlik anomalilerin tespitinde zaman serilerinin orjinal durumları kullanılırken, şekil anomalilerinin tespitinde ise zaman serilerinin otokorelasyon gösterimleri kullanılmıştır.

25

1 Araştırmacılar sabit uzunlukta bir aralık penceresi
2 ile bir alt dizi seti oluşturmuş ve bu alt diziler
3 içerisindeki normal yapıyı göstermek için **FCM**
4 algoritmasından yararlanmıştır. Her bir alt dizinin
5 bir kümeye benzersizliğini ölçmek için yeniden
6 yapılandırma kriteri kullanmışlar ve yeniden
7 yapılandırma hatasını hesaplayarak anomali skoruna
8 atamışlardır.

11 2.2 Makale 2: “Large-Scale Unusual Time 12 Series Detection”[6]

13 Bu makalede araştırmacılar büyük ölçekte ki zaman
14 serileri kullanarak anomali tespiti üzerine çalışma
15 yapmışlardır. Bu çalışmada zaman verileri binlerce
16 sunucunun saat başı toplanan performans
17 ölçümlerinden oluşmaktadır. Ve araştırmacılar da
18 olağan dışı davranan sunucuları tespit etmeyi
19 hedeflemektedirler.

20 Araştırmacılar, her zaman serisi için karakteristik
21 ölçümleri hesapladıktan sonra öznitelik vektörlerini
22 çıkartmışlardır. Elde edilen öznitelik vektörleri
23 üzerinde ana bileşen ayrışımı uygulayarak, çeşitli iki
24 değişkenli anomali tespiti yöntemlerini
25 kullanmışlardır. Bu yaklaşım, en olağan dışı zaman
26 serilerini öznitelik vektörleri üzerinden
27 belirlemelerini sağlamıştır. Yazarlar, kullandıkları
28 yöntemlerin en yüksek yoğunluklu bölgelere ve α -
29 hull'lara dayandığını ifade etmektedirler.

30 Araştırmacılar bu yazıda büyük miktardaki zaman
31 serilerinde olağandışı zaman serilerini tespit etmek
32 için Ana Bileşen Analizi (Principal Component
33 Analysis) ile beraber çok-boyutlu anomali tespiti
34 kullanmayı önermektedirler. Bu yaklaşımlarının
35 80% isabet oranına sahip olmakla beraber, 1000
36 zaman serisini işlemek için 0.5 saniyeden daha azına
37 ihtiyaç duyduğunu belirterek bu yöntemlerinin güçlü
38 ve isabetli olduğunu vurgulamaktadırlar.

42 2.3 Makale 3: “Anomaly Detection of 43 Spacecraft Based on Least Squares Support 44 Vector Machine”[7]

45 Bu yazı, yörüngede bulunan uzay araçlarında ortaya
46 çıkan anomalileri tespit etmek için **Least Squares**
47 **Support Vector Machine (LS-SVM)**
48 algoritmasına dayanan bir yöntem sunmaktadır.
49 Voltaj, sıcaklık, akım, titreşim gibi veriler toplanıp
50 ön işleme tabii tutulduktan sonra öznitelikleri
51 çıkartılıp Ana Bileşen Analizi kullanılarak öznitelik
52 dizileri oluşturulmuş ve LS-SVM kullanılarak uzay
53 aracındaki anormal davranışlar tespit edilmeye
54 çalışılmıştır.

56

57 Bu çalışma, uzay aracındaki anormal durumları
58 tespit etmek amacıyla Least Squares Support Vector
59 Machine (LS-SVM) tabanlı bir yöntem
60 önermektedir. PCA tabanlı öznitelik çıkarma ve LS-
61 SVM kullanılarak gerçekleştirilen anormal durum
62 tespiti sonuçları, yüksek doğruluk ve etkinlik
63 göstermektedir. Ancak, yöntemin sınırlamaları ve
64 PCA tabanlı öznitelik çıkarma yöntemlerinin
65 geliştirilmesi ihtiyacı vardır.

67 2.4 Makale 4: “Probabilistic anomaly 68 detection in natural gas time series data”[8]

69 Bu makale doğalgaz zaman serisi verileri için bir
70 olasılıksal yaklaşım sunmaktadır. Bir dizi
71 anomalinin sebepleri incelenip kategorize edilmiş ve
72 Bayesian maksimum olabilirlik sınıflandırıcısına
73 bilinen anomalilerin geçici yapısı öğretilmiştir.
74 Bilinmeyen bir zaman serisi verisi sunulduğunda,
75 sistem, hava durumu girdilerini kullanarak bir lineer
76 regresyon modeli ile anormallikleri tespit
77 etmektedir. . Daha sonra, anormallikler yanlış
78 pozitiflerle test edilir ve Bayesian bir sınıflandırıcı
79 kullanılarak sınıflandırılır. Bu yöntem aynı zamanda
80 bilinmeyen kaynaklı anormallikleri tanımlayabilir.
81 Bu nedenle, bir veri noktasının anormal olma
82 olasılığı, hem bilinen hem de bilinmeyen kaynaklı
83 anormallikler için sağlanmaktadır.

84 Bu yöntem günlük tüketim verilerine
85 uygulandığında, veri temizleme yoluyla kök
86 ortalama kare hatalar (RMSE'ler) için %37.5 ve
87 ortalama mutlak yüzde hataları (MAPE'ler) için
88 %7.84'ün üzerinde ortalama bir iyileşme sağlar.
89 Çalışma, veri imputasyonu için güçlü tahmin
90 modellerinin kullanılması ve Bayesian
91 sınıflandırıcıyı geliştirmek için dış kaynaklı
92 girdilerin eklenmesi gibi daha fazla iyileştirmeler
93 önermektedir. Yöntem, zaman serisi verilerinde
94 bilinen dış kaynaklı faktörlere sahip çeşitli alanlarda
95 uygulanabilirliğini göstermektedir.

97 2.5 Makale 5: “Generic and Scalable 98 Framework for Automated Time-series 99 Anomaly Detection”[9]

100 Bu makalede Yahoo'nun kendi kullandığı, büyük
101 ölçekli zaman serili veriler için açık kaynaklı bir
102 otonom anomali tespit kütüphanesi tanıtılmaktadır.
103 Bu kütüphane (**EGADS**) mevcut çoğu anomali tespit
104 yaklaşımlarının büyük sorun yaşadığı
105 ölçeklenebilirlik konusunda birden fazla tespit ve
106 tahmin modelleri kullanarak bir çözüm ürettiğini
107 belirtmektedir. Gerçek ve yapay veri kullanarak

kendi yaklaşımlarının diğer sistemlerden 50-60% oranında daha iyi olduğunu keşfetmişlerdir.

2.6 Makale 6: “ANOMALY DETECTION IN AIRCRAFT DATA USING RECURRENT NEURAL NETWORKS (RNN)”[10]

Bu makale, uçakların Uçuş Veri Kaydedici (FDR) veya Uçuş Operasyonel Kalite Güvence (FOQA) verilerinden toplanan çoklu değişkenli, zaman serisi verilerinde anomalilerin tespiti üzerine odaklanmaktadır. Endüstri standardı "Aşma Algılama" algoritması, belirli parametrelerin ve eşiklerinin listesini kullanarak bilinen sapmaları tanımlar. Buna karşılık, Makine Öğrenimi algoritmaları, verideki bilinmeyen olağandışı desenleri yarı gözetimli veya gözetimsiz öğrenme yoluyla tespit eder.

Makale, Recurrent Neural Networks (RNN)'in Long Term Short Term Memory (LTSM) ve Gated Recurrent Units (GRU) mimarileriyle uygulandığını tanımlar. RNN algoritmalarının hassasiyet=1, anımsama=0.818 ve F1 skoru=0.89 değerleriyle 11 anomaliden 9 tanesinin tespit edildiği belirtilmektedir.

2.7 Makale 7: “Hybrid robust convolutional autoencoder for unsupervised anomaly detection of machine tools under noises”[11]

Bu makale gürültülü ortamlarda makinelerden etiketli veri elde etmenin zorluğuna bir çözüm üretmek için gürültülü içerisinde çalışan makinelerin denetimsiz anomali tespiti için hibrit sağlam bir konvolüsyonel otokodlayıcı (HRCAE) geliştirilmesinden bahsetmektedir. Bu bahsedilen method bir Paralel Konvolüsyonel dağılım Uydurma (PCDF) modülü oluşturur ve bir . Birleştirilmiş Yönlü Mesafe (FDD) kayıp fonksiyonu tasarlayarak anomali tespitinin sağlamlığını artırır. Diğer denetimsiz otokodlayıcı yöntemleriyle karşılaştırıldığında bu metodun farklı gürültülerin arasında daha iyi tespit yapabildiği belirtilmiştir.

2.8 Makale 8: “Vibration-based anomaly detection using LSTM/SVM approaches”[12]

Bu çalışma, yarı-üstünlü anomali tespiti için yeni veri odaklı mimariler önermektedir. LSTM regresörleri ve tek sınıf SVM sınıflayıcılarını

birleştiren bu mimariler, dişli aşınması ve rulman arızalarının titreşim sinyallerini hedeflemektedir. Standart bir model, dişli arızalarını tespit etmede etkili olsa da, rulman arızalarının neden olduğu zayıf değişiklikleri tespit etmekte zorlanır. Bu zorluğu aşmak için, sağlıklı sinyallerle eğitilmiş iki aşamalı bir model önerilmektedir.

Bu model, sinyalin sağlıklı bileşenlerini kaldıran birinci LSTM regresörü ve ardından artık sinyaldeki belirgin bileşenleri temizleyen ikinci bir LSTM regresörü içerir. Bu ekleme, rulman arızasının ilerlemesiyle sinyalin rastgele bileşenlerindeki bir artışı tanıma yeteneği sağlar. Sonuç olarak, sürekli olmayan veri kümeleri için ise, otoregresif bir LSTM eklemenin avantajı olmadığı ve basit istatistiksel özelliklere dayalı tek sınıf SVM aykırı değer tespitin daha etkili ve güvenilir olduğu gösterilmiştir.

2.9 Makale 9: “A Deep Support Vector Data Description Method for Anomaly Detection in Helicopters ”[13]

Bu makalede helikopterlerde anomali tespiti için Support Vector Data Description (SVDD) yöntemi baz alınmıştır. Bu çalışmada da otonom bir anomali tespiti kurmak için derin bir SVDD modeli sunulmuştur. SVDD modeline uygun öznetelik çıkarmak içinde bir CNN uygulanmıştır.

Deneyler, bu çalışmada kullandığı AirbusSAS tarafından sağlanan helikopter titreşim veri seti üzerinde yürütülmüş ve F1 skorunun 94%'e kadar ulaştığını göstererek yöntemin etkili olduğunu göstermiştir. Karşılaştırmalı analiz, önerilen yöntemin diğer giriş tiplerinden daha iyi performans gösterdiğini ortaya koymuştur.

2.10 Makale 10: “ANOMALY DETECTION AND CLASSIFICATION IN TIME SERIES WITH KERVOLUTIONAL NEURAL NETWORKS ”[14]

Bu makale, kervolutional sinir ağlarının (KCNN'ler) zaman serisi verilerine uygulanan potansiyelini araştırıyor. KCNN'ler, konvolüsyonel sinir ağlarının (CNN'ler) uzun mesafeli ilişkileri yakalama yeteneklerini, CNN'lerin yerel özelliklerini yakalama yetenekleriyle birleştiren yeni bir tür koniyeldir.

Makale, KCNN'lerin CNN'lerden daha iyi performans gösterdiğini gösteren iki deney rapor ediyor. İlk deneyde, KCNN'ler zaman serilerinde bir sınıflandırma görevinde test edilmiş ve CNN'lerden

1 daha yüksek doğruluk elde etmişlerdir. İkinci
2 deneyde ise KCCNN'ler helikopterlerde ivme ölçerler
3 tarafından kaydedilen zaman serisi verilerinde
4 anomalileri tespit etmek için test edilmiş ve
5 CNN'lerden daha iyi performans göstermişlerdir. Bu
6 sonuçlar, KCCNN'lerin zaman serisi verilerinde arıza
7 tespiti ve sınıflandırma için etkili bir araç
8 olabileceğini göstermektedir..

9 10 **2.11 Makale 11: “Exemplar Learning for** 11 **Extremely Efficient Anomaly Detection in** 12 **Real-Valued Time Series ”[15]**

13
14 Bu makale, gerçek değerli zaman serilerindeki
15 anomalileri tespit etmek için yeni bir algoritma
16 sunmaktadır. Algoritma, bir eğitim zaman serisinde
17 bulunan çeşitli alt dizileri temsil etmek için küçük
18 bir örnek seçimine dayanır. Örnekler, benzer alt
19 dizilerde bulunan karakteristik yörüngeyi (düşük
20 frekans bilgisi) ve stokastik özelliklerini (yüksek
21 frekans bilgisi) yakalayan İstatistiksel ve Pürüzsüz
22 Yörünge özelliklerini kullanır. Makale, yeni
23 algoritmanın önceki algoritmalarından çok daha hızlı
24 ve daha doğru olduğunu göstermektedir.

25 26 **2.12 Makale 12: “Detecting Anomalies in a** 27 **Time Series Database”[16]**

28
29 Bu makalede yazarlar yarı denetimli anomali tespit
30 tekniklerini dair yaptıkları değerlendirmeleri
31 sunmuşlardır. Bu değerlendirmenin sonunda şu
32 sonuçlara varmışlardır;

- 33 • Sürekli zaman serilerinde çalışan teknikler,
34 genellikle ayrık diziler üzerinde çalışan
35 tekniklerden üstündür.
- 36 • Model bağımsız olan çekirdek ve pencere
37 tabanlı teknikler, zaman serisi verileri için
38 bir model oluşturmaya çalışan tahmini ve
39 segmentasyon tabanlı tekniklerden daha iyi
40 performans gösterme eğilimindedir.
- 41 • Çekirdek tabanlı teknikler, pencere tabanlı
42 tekniklere göre daha hızlıdır.
- 43 • Normal ve anormal verilerin doğası ile farklı
44 tekniklerin performansı arasında birkaç
45 ilişki vardır.

46 47 **2.13 Makale 13: “Anomaly Detection** 48 **Method of Aircraft System using** 49 **Multivariate Time Series Clustering and** 50 **Classification Techniques”[17]**

51
52 Bu makalede uçak sistemlerinde sapma tespiti için
53 yeni bir yöntem önerilmiştir. Yöntem, uçuş evreleri
54 içindeki değişkenlerin döngüselliklerini dikkate alarak
55 normal ve anormal davranışları ayırt eder. Yöntem,

56 DBSCAN ve DTW algoritmalarını kullanarak
57 yüksek doğrulukla çalışır.

58
59 Önerilen yöntemde, bu anormal ve normal
60 davranışlar karşılaştırıldığında, yeni bir önem puanı
61 kullanılarak açıklanır. Davranışlar, Zaman Serisi
62 Ormanı (TSF) algoritması ve silüet kriteri
63 kullanılarak belirlenir. Metod, Bombardier'in Uçak
64 Sağlık İzleme Sisteminden bir örnek kullanılarak
65 eğitilmiş ve test edilmiştir. Normal ve anormal
66 davranışları kümelenme silüet puanı 0,95 olacak
67 şekilde ayırt eder ve bilinmeyen davranışları %89
68 hassasiyetle tespit eder.

69 70 **2.14 Makale 14: “Anomaly Detection** 71 **Method of Aircraft System using** 72 **Multivariate Time Series Clustering and** 73 **Classification Techniques”[18]**

74
75 Bu makalede gerçek dünyadaki verilerde anomaliler
76 seyrek olarak gerçekleştiğinden, çok değişkenli
77 zaman serilerinde sapma tespitinin zorluğu ele
78 alınmıştır. Yazarlar bir seferde yalnızca bir özellik
79 kümesini dikkate alan bir özellik toplama tekniği
80 önermektedir. Bu, anomalilerin daha küçük bir
81 özellik kümesinde daha olası olduğundan, anomali
82 tespiti için daha etkilidir.

83
84 Ek olarak, yaklaşımın etkinliğini ve
85 genelleştirilmesini iyileştirmek için, yazarlar PCA
86 ile hesaplanan iç içe geçmiş rotasyona dayalı bir
87 dönüşüm uygularlar. Bu, anomalileri daha iyi ayırt
88 etmeye yardımcı olur. Son olarak, temel modellerin
89 çıktılarını birleştirmek için, yazarlar bir Lojistik
90 Regresör kullanan yarı denetimli bir yaklaşım
91 önerirler. Bu, tahmin performansını daha da artırır.
92 Genel olarak, önerilen yöntem, çok değişkenli
93 zaman serilerinde sapma tespiti için etkili bir
94 yaklaşımdır.

95 96 **2.15 Makale 15: “Operational Anomaly** 97 **Detection in Flight Data Using a Multivariate** 98 **Gaussian Mixture Model”[19]**

99
100 Bu makale, helikopter uçuş verilerinde sapmaları
101 tespit etmek için GMM'leri kullanma yöntemini
102 önermektedir. Önerilen yöntem, öncelikle veri
103 kümesi içindeki özelliklerin sayısını azaltmak için
104 bir boyutsallık azaltma tekniği kullanır. GMM daha
105 sonra azaltılmış boyutlu veri kümesi üzerinde
106 eğitilir. Eğitilmiş GMM daha sonra yeni veri
107 noktalarını puanlamak için kullanılır ve düşük puan
108 alan veri noktaları anomali olarak işaretlenir.

1 Önerilen yöntem, Maryland Üniversitesi'nden
2 helikopter uçuş verisi kümesi üzerinde
3 değerlendirilmiştir. Sonuçlar, yöntemin yüksek
4 doğrulukla anomalileri tespit edebildiğini
5 göstermektedir. Yöntem ayrıca hesaplama açısından
6 verimlidir, bu da onu gerçek zamanlı uygulamalar
7 için uygun hale getirir. Genel olarak, bu metodun
8 helikopter uçuş verilerinde sapma tespiti için umut
9 verici bir yaklaşım olduğu belirtilmektedir.
10 Helikopter operasyonlarının güvenliğini ve
11 güvenilirliğini iyileştirme potansiyeline sahiptir.

13 2.16 Makale 16: “Hybrid approach for 14 Anomaly Detection in Time Series Data”[20]

16 Bu makalede, zaman serisi verilerinde anomali
17 tespiti için bir hibrit yaklaşım önerilmektedir.
18 Önerilen yaklaşım, LSTM Autoencoder ve doğrusal
19 sınıflandırıcıyı birleştirir. LSTM Autoencoder,
20 zaman serisinden verimli temsiller öğrenir. Bu
21 temsiller, doğrusal sınıflandırıcı tarafından
22 anomalileri tespit etmek için kullanılır.

24 Deneyler, gerçek dünya veri kümelerinde yapılan
25 karşılaştırmalara göre, önerilen hibrit yaklaşımın
26 mevcut en iyi yöntemlerden daha iyi performans
27 gösterdiğini göstermektedir.

29 2.17 Makale 17: “A Method to Handle 30 Unstable Time Series in Anomaly Detection 31 Problem”[21]

33 Bu makalede, zaman serisi varyasyonunu azaltmak
34 için bir yöntem önerilmektedir. Bu yöntem, zaman
35 serisinden kararsız parçaları otomatik olarak arayan
36 ve hariç tutan bir algoritmadan oluşur. Algoritma,
37 ham veri örneklerinin toplanmasından makine
38 öğrenimi modelleriyle gelecekteki çalışmalar için
39 vektörlemelerine kadar adımları kapsayan bir ön
40 işleme algoritmasının parçası olarak
41 uygulanmaktadır.

43 Önerilen yöntem, gaz sensörlerinden gelen zaman
44 serilerinde test edilmiştir. Bu zaman serileri,
45 sensörün hassas yüzeyinin çevrede bulunan çeşitli
46 gazlara karşı reaksiyonunu yansıtmaktadır.
47 Deneyler, yöntemin mevcut en iyi yöntemlerden
48 daha iyi performans gösterdiğini göstermektedir.

50 2.18 Makale 18: “Anomaly Detection for Key 51 Performance Indicators Through Machine 52 Learning”[22]

54 Bu makale ağ hizmetlerinin stabil olup olmadığını
55 anlamak için çeşitli kilit performans göstergeleri

56 üzerinden anomali tespiti üzerinde çalışma
57 yapmıştır. Sorun iki kısma ayrılmıştır. İlk kısımda
58 zaman serisinde bir sonraki noktayı tahmin etmek
59 için zaman serisi analiz metodu olan **Holt-Winters**,
60 **ARIMA** modeli, regresyon tabanlı **GBRT** ve
61 **LSTM** teknikleri kullanılmıştır. Sonrasında ise
62 anomali kuralı belirlenmektedir. Ve tahmin edilen
63 değer ile asıl değer karşılaştırılarak asıl değer
64 anomali olup olmadığı tespit edilir.

66 Bahsedilen metodlar arasında en iyi performansa
67 sahip olan tekniğin, en düşük **MSE**, en yüksek **F1**
68 skoru ve en kısa eğitim süresine sahip olan **GBRT**
69 olduğu belirtilmektedir.

71 2.19 Makale 19: “A Comparative Study of 72 Cluster Based Outlier Detection, Distance 73 Based Outlier Detection and Density Based 74 Outlier Detection Techniques ”[23]

76 Bu makalede yazarlar kümeleme bazlı, uzaklık bazlı
77 ve yoğunluk bazlı anomali tespit teknikleri üzerine
78 bir karşılaştırma araştırması yapmaktadırlar.
79 Araştırmacılar bu tekniklerin performanslarını
80 ölçmek için farklı türde gürültü ve anomali içeren
81 çeşitli veri setleri kullanmaktadırlar. Her bir tekniğin
82 kendine has parametreleri olduğundan direkt bir
83 karşılaştırmanın mümkün olmadığı sonucunu
84 çıkarmışlardır.

86 2.20 Makale 20: “Hybrid Machine Learning 87 for Anomaly Detection in Industrial Time- 88 Series Measurement Data”[24]

90 Bu makale, çoklu değişkenli zaman serisi ölçüm
91 verilerinden oluşan büyük veri kümelerinde
92 anomalilerin tespiti için yeni bir yöntem
93 önermektedir. Yöntem, anahtar performans
94 göstergeleri (KPI'lar) ve uzun kısa süreli bellek
95 (LSTM-VAE) varyasyonel otoenkoder adlı iki farklı
96 makine öğrenimi tekniğini birleştirir. KPI'lar, sistem
97 davranışını tahmin etmek için kullanılır. **LSTM-VAE** ise verilerin istatistiksel belirsizliğini dikkate
98 alarak sistemin normal davranışını modeller. Bu iki
99 yöntemin kombinasyonu, verilerde anomalileri daha
100 güvenilir bir şekilde tespit etmeyi sağlar.

102 Yöntem, bina temelleri için bir zemin iyileştirme
103 işlemiyle ilgili bir vaka çalışması ile doğrulanır.
104 Sonuçlar, yöntemin verilerde anomalileri tespit
105 etmede etkili olduğunu göstermektedir.

3. YÖNTEM

Helikopter verilerinde anomali tespiti için beş adet model eğitilmiştir. Modeller eğitilmeden önce veri ön işlemeye tâbi tutulmuştur.

3.1 Ön İşleme

Eğitim veri seti doğası gereği yüksek boyutludur. Eğitim verisi 1677 satır 61440 kolondan, test verisi ise 594 satır 61440 kolondan oluşmaktadır. Ayrıca test verisi için 594 satır 2 kolonluk etiket verisi bulunmaktadır. Etiketler için 0 normal, 1 anomaliyi ifade etmektedir.

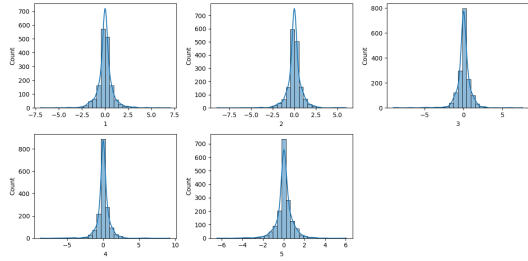
Tablo 1. Veri setinin Tanımı

<i>RangeIndex: 1677</i>	<i>RangeIndex: 594</i>
<i>entries, 0 to 1676</i>	<i>entries, 0 to 593</i>
<i>Columns: 61440</i>	<i>Columns: 61440</i>
<i>entries, 0 to 61439</i>	<i>entries, 0 to 61439</i>
<i>dtypes: float64(61440)</i>	<i>dtypes: float64(61440)</i>

Eğitim Seti

Test Seti

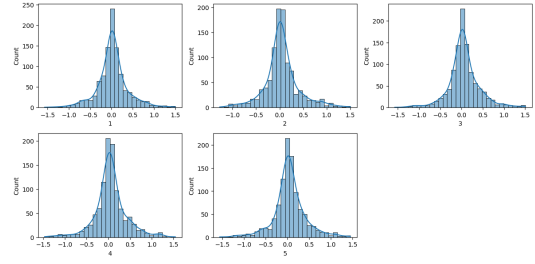
3.1.1 Aykırı Değerleri Silme



Şekil 1. Eğitim verisinin düzensiz hali

Şekil 1'de eğitim verisinin ön işlemeye tâbi tutulmadan önceki hali sunulmuştur. Bu grafikte ilk 5 ögenin dağılımını gösterilmektedir. Grafikte değerlerin dağılımının oldukça düzensiz olduğu görülmektedir. Ve -7.5 ile 7.5 arasında değişmektedir. Ayrıca değerler arasındaki korelasyon da oldukça zayıftır. Bu bilgiler doğrultusunda verinin heterojen olduğu çıkarımı yapılabilir.

Veri setinde mevcut bulunan gürültüyü ve düzensizliği ortadan kaldırmak için verilerin Z-Skoru hesaplanmıştır. Bu skordan büyük olan veriler düzensizliğe sebep olmaktadır ve modellerin iyi bir şekilde eğitilmesi için veri setinden çıkartılmıştır.



Şekil 2. Eğitim verisinin düzenli hali

Z-Skoru 3'ten büyük olan veriler silindikten sonra Şekil 2'de görüldüğü üzere veri seti daha düzenli ve homojen hale gelmiştir. Satır sayısı 1038'e düşmüştür.

3.1.2 Eksik değerlerin doldurulması

Veri ön işlemenin önemli adımlarından biri de eksik değerlerin imputasyonudur. Bu çalışma da üç adet imputasyon yöntemi kullanılmıştır. İlk olarak Ortalama İmputasyonu kullanılmıştır. Bu metodla eksik değerler, değişkenin ortalama değeri ile doldurulmuştur. Tablo 2 ve 3'te ilk 5 sütun için ortalama imputasyondan öncesi ve sonrası paylaşılmıştır.

	0	1	2	3	4
count	1038	1038	1038	1038	1038
mean	0.0034	0.0224	0.0465	0.0577	0.0548
std	0.3618	0.3641	0.3676	0.3715	0.376
min	-1.553	-1.568	-1.273	-1.574	-1.47
25%	-0.147	-0.121	-0.101	-0.11	-0.0968
50%	0.017	0.0187	0.0252	0.0362	0.036
75%	0.159	0.1645	0.1824	0.2073	0.1969
max	1.407	1.4994	1.5049	1.5083	1.5432

Tablo 2. Ortalama İmputasyonu öncesi

	0	1	2	3	4
count	1038	1038	1038	1038	1038
mean	0.003397	0.022433	0.046546	0.057736	0.054794
std	0.361824	0.364087	0.367609	0.371529	0.375983
min	-1.553164	-1.567793	-1.272794	-1.574047	-1.470111
25%	-0.1479	-0.121382	-0.100552	-0.110418	-0.096798
50%	0.016996	0.018709	0.025158	0.036214	0.035987
75%	0.159156	0.164466	0.182395	0.207288	0.196914
max	1.406972	1.499392	1.504855	1.508307	1.54317

Tablo 3. Ortalama İmputasyonu sonrası

Yukarıdaki tablolarda sütun isimleri şöyle açıklanır:

- **count:** satır sayısı,
- **mean:** her sütundaki değerlerin ortalaması,
- **std:** her sütundaki değerlerin standart sapması,
- **min:** en küçük değer,
- **25%:** her sütundaki değerlerin alt çeyrek yüzdeliği,
- **50%:** her sütundaki değerlerin medyanı,

- **75:** her sütundaki değerlerin üst çeyrek yüzdeliği,
- **max:** en büyük değer.

Tablo 2 ve 3 incelendiğinde ortalama imputasyonunda kayda değer bir değişiklik olmadığı görülmektedir. Bu da ortalama imputasyonunu veri setine etkisinin sınırlı olduğu anlamına gelmektedir.

İkinci yöntem olarak $n_neighbors = 5$ parametresiyle **K-NN** imputasyonu kullanılmıştır. Bu imputasyon sonrası sonuçlar Tablo 4'te paylaşılmıştır.

	0	1	2	3	4
count	1038	1038	1038	1038	1038
mean	0.003397	0.003397	0.046546	0.057736	0.054794
std	0.361824	0.361824	0.367609	0.371529	0.375983
min	-1.553164	-1.553164	-1.272794	-1.574047	-1.470111
25%	-0.1479	-0.1479	-0.100552	-0.110418	-0.096798
50%	0.016996	0.016996	0.025158	0.036214	0.035987
75%	0.159156	0.159156	0.182395	0.207288	0.196914
max	1.406972	1.406972	1.504855	1.508307	1.54317

Tablo 4. K-NN İmputasyonu sonrası

	0	1	2	3	4
count	1038	1038	1038	1038	1038
mean	0.003397	0.022433	0.046546	0.057736	0.054794
std	0.361824	0.364087	0.367609	0.371529	0.375983
min	-1.553164	-1.567793	-1.272794	-1.574047	-1.470111
25%	-0.1479	-0.121382	-0.100552	-0.110418	-0.096798
50%	0.016996	0.018709	0.025158	0.036214	0.035987
75%	0.159156	0.164466	0.182395	0.207288	0.196914
max	1.406972	1.499392	1.504855	1.508307	1.54317

Tablo 5. F-B Fill İmputasyonu sonrası

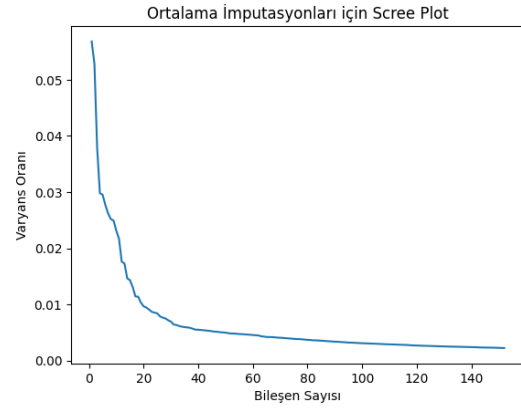
Tablo 2, Tablo 4 ve Tablo 5 karşılaştırıldığında ilk 5 sütun için yine büyük bir değişiklik söz konusu değildir. Bu durum, imputasyon yöntemlerinin belirli bir istatistiksel özet üzerinde benzer etkilere sahip olduğunu düşündürülebilir. Ancak, bu değerlerin benzer olması, imputasyon yöntemlerinin tümüyle aynı sonuçları ürettiği anlamına gelmez.

3.1.2 Öznitelik Çıkarımı

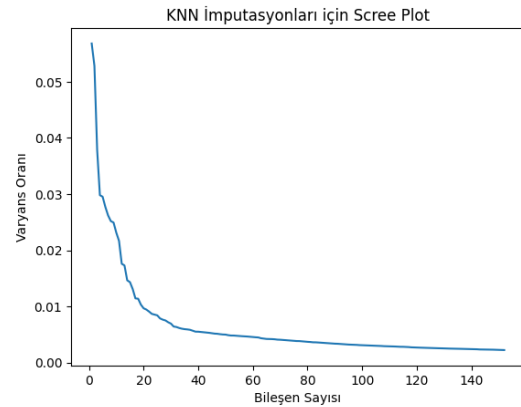
Makine öğrenmesinde öznitelik, bir fenomenin ölçülebilir özelliği ya da karakteristiğidir [25]. Bu veri setinde yapısı gereği yüksek boyutta öznitelik bulunmaktadır. Aşırı öznitelik sayısı, modelin genelleme performansını düşürebilir ve aşırı öğrenmeye neden olabilir. İdeal öznitelik sayısı, hem yeterli bilgi sağlamalı hem de modelin karmaşıklığını kontrol altında tutmalıdır. Buna istinaden tek tek öznitelik çıkarılması bu veri seti için mümkün olmadığından yaygın bir öznitelik çıkarma metodu olan **PCA** kullanılmıştır.

PCA'nın uygulanması için özniteliklerden oransal olarak ne kadar çıkarım yapılacağı belirtilmelidir. Burada dikkat edilmesi gereken nokta az öznitelikle verisetinin çoğunluğunu örnekleyecek şekilde bir oran seçilmesidir. Veri setinin ne durumda olduğunu görmek için ilk olarak standart bir değer olan $n_components=0.7$ parametresiyle PCA uygulanmıştır. PCA uygulanmadan önce eğitim ve test setleri (1038,61440) ve (594,61440) şekline sahiptir.

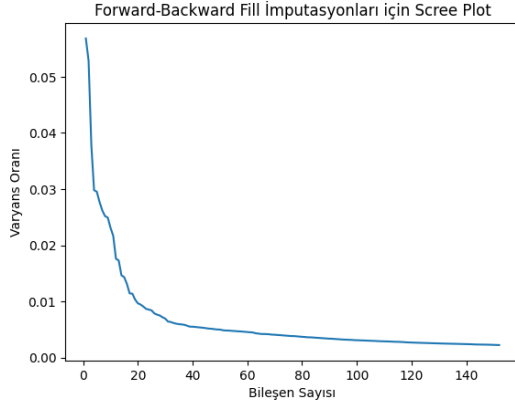
PCA uygulandıktan sonra ise eğitim ve test setleri (1038,152) ve (594,152) şekillerine sahip olmuşlardır. Ayrıca PCA'nın uygulanması için veri setlerine Min-Max yöntemiyle normalizasyon uygulanmıştır.



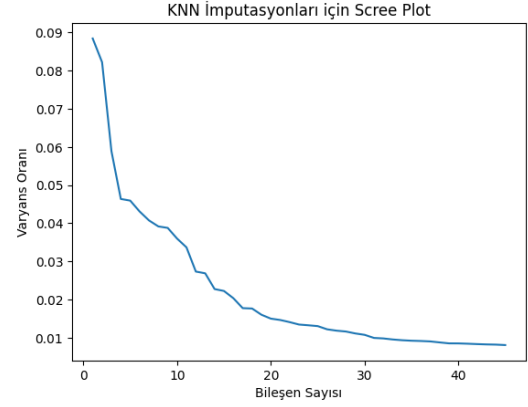
Şekil 3.



Şekil 4.



Şekil 5.

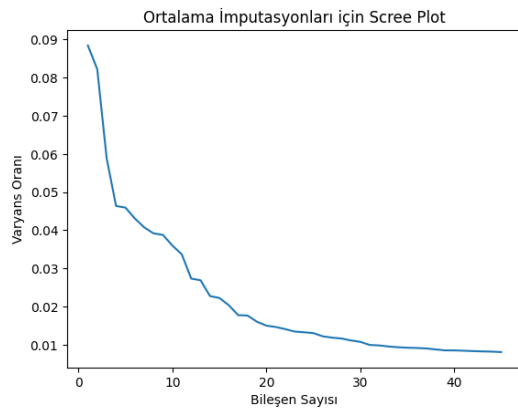


Şekil 7.

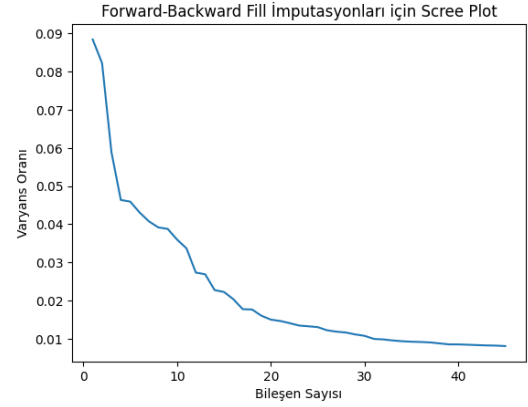
Şekil 3-5'te görülen grafiklerde eksenlerin açıklaması aşağıda verilmiştir;

- **bileşen:** Analiz sonrasında elde edilen bileşen sayısını ifade eder. PCA, veri setindeki değişken sayısını azaltmak için kullanılır, ve bu bileşen sayısı analizin sonucunda belirlenir.
- **varyans:** Her bir bileşenin toplam varyans içindeki oranını gösterir. Yüksek varyans açıklama oranına sahip olan bileşenler, orijinal veri setindeki önemli bilgileri korur.

Şekil 3-5'te görüldüğü gibi varyans oranı 40. bileşenden sonra azalmaya başlamıştır. Bu ilk 40 bileşenin sahip olduğu bilgi oranının, veri setinin genelini kapsadığı anlamına gelmektedir. Çok düşük varyans oranına sahip bileşenleri modele beslemek aşırı uyum riskini ortaya çıkarır. Bu yüzden PCA 0.45 bileşen oranı ile tekrar uygulanmıştır.



Şekil 6.



Şekil 8.

Şekil 6-8'te görüldüğü üzere bileşen oranı 0.45 yapıldığında Şekil 3-5'te görülen düşük varyanslı bileşenler eğitim ve test setlerinden çıkarılmıştır. Öznitelik çıkarımı yapılarak eğitim seti model eğitmeye uygun hale gelmiştir.

3.2 Model Eğitimi

Bu çalışmada dört adet makine öğrenmesi modeli ve bir adet derin öğrenme modeli, üç farklı şekilde impute edilmiş verilerle eğitilmiştir. Bu modellerin karşılaştırılmasında **Kesinlik (A)**, **Hassasiyet (P)**, **Duyarlılık (R)** ve **F₁ Skoru** metrikleri kullanılmıştır. Ayrıca derin öğrenme modeli için **Receiver Operator Characteristics (ROC)** eğrisi de kullanılmıştır.

3.2.1 Isolation Forest

Isolation Forest, ilk olarak 2008'de Fei Tony Liu tarafından anomali tespiti için geliştirilmiştir [26]. Bu algoritma anomali tespitini ikili ağaçlar kullanarak tespit etmektedir [27][28].

Bu model eğitilirken yapılan ince ayar sonucunda en iyi hiper parametreler şu şekilde açıklanmıştır;

- **n_estimators=200:** Toplamda oluşturulan izolasyon ağaçlarının sayısını belirler.
- **max_samples=0.2:** Her bir izolasyon ağacının oluşturulurken kullanılacak örnek sayısını belirler.
- **contamination=0.175:** Veri kümesindeki anomali oranını belirler. Bu oran, izolasyon ağaçları tarafından izole edilmiş örneklerin oranını ifade eder.
- **random_state=42:** Modelin tekrarlanabilir olmasını aynı parametrelerle çalıştırıldığında aynı sonuçların elde edilmesini sağlar.

3.2.2 One-Class SVM

Bu algoritmanın çalışma prensibi, önce eğitim verilerinde yalnızca normal sınıfa ait olan örnekleri kullanarak bir model oluşturmaktır. Bu model, normal sınıfa ait örnekleri içeren bir bölge veya hacim tanımlar. Daha sonra, bu bölge dışında kalan noktaları aykırı değer olarak kabul eder.

Bu bölgeyi oluşturmak için seçilen yapılan ince ayarda en iyi parametreler şunlar seçilmiştir;

- **nu=0.11:** Modelin eğitimi sırasında hata toleransını kontrol eder.
- **tol=0.008:** Eğitim süreci bu tolerans eşiği altına düştüğünde durur.
- **kernel=rbf:** SVM'in veri kümesindeki örnekleri daha yüksek boyutlu bir özellik uzayına haritalamak için kullanılan bir matematiksel işlemdir. (Radial Basis Function)
- **gamma_value=scale:** RBF çekirdeğinin genişliğini kontrol eden bir parametredir. "scale" seçeneği, gamma'nın $1 / (n_feature * X.var())$ şeklinde ölçeklendiği anlamına gelir.

3.2.3 K-Means Cluster

Bu algoritma, veri noktalarını belirli sayıda küme içinde gruplandırmak için kullanılan bir kümeleme algoritmasıdır. Temel olarak, veri noktalarını birbirine benzerliklerine göre gruplandırmaya çalışır. Belirli sayıda küme merkezi rastgele seçilir. Her veri noktası en yakın küme merkezine atanır. Küme merkezleri, kendi içindeki veri noktalarının ortalaması olarak güncellenir. Atama ve güncelleme işlemleri, küme merkezleri değişmeyene veya belirli bir iterasyon sayısına ulaşana kadar tekrarlanır [29].

Bu algoritma ile eğitilen modelin parametreleri şu şekilde açıklanmıştır;

- **n_cluster=6:** Küme sayısını ifade eder.
- **n_init=8:** Bu algoritma rastgele başlangıç merkezleri seçer ve her bir başlangıçta algoritmayı çalıştırır. *n_init* parametresi, bu rastgele başlangıçların kaç kez deneneceğini belirtir ve en iyi sonucu seçer.
- **algorithm=lloyd:** K-Means algoritmasının kullanacağı algoritmayı belirler. 'lloyd', K-Means'in klasik Lloyd's algoritması olduğunu belirtir.
- **threshold=99:** Anomali tespiti için eşik değerini ifade eder.

3.2.4 Gaussian Mixture Model

Gaussian Mixture Model (GMM), istatistiksel bir modelleme tekniğidir ve veri kümesini birden fazla Gauss (normal) dağılımının bir kombinasyonu olarak temsil eder. Her bir Gauss bileşeni, veri noktalarının belirli bir kısmını temsil eder ve ağırlıkları, merkezi konumları ve kovaryans matrisleri kullanılarak tanımlanır. GMM, esnek bir olasılık dağılımı sağlar ve karmaşık veri yapılarını modellemek, veri setindeki farklı grupları tanımlamak için kullanılır [30]. Bu modelin eğitilmesinde kullanılan parametreler şunlardır:

- **n_components=110:** bileşen sayısı,
- **covariance_type=full:** kovaryans tipi,
- **threshold=40:** eşik değeri.

3.2.5 LSTM

Long Short-Term Memory (LSTM), özellikle zaman serisi verileri gibi uzun vadeli bağımlılıkları işlemek için tasarlanmış bir tür rekürrent sinir ağı (RNN) yapısıdır. LSTM, geleneksel RNN'lerin karşılaştığı uzun vadeli bağımlılık sorununu çözmek için özel hücre yapıları kullanır. Bu hücreler, giriş verisinin anlık değeri yanı sıra önceki zaman adımlarından gelen bilgileri koruyarak ve güncelleyerek çalışırlar. Bu özellikleri sayesinde LSTM, uzun vadeli bağımlılıkları etkili bir şekilde modelleyebilir ve bu nedenle özellikle zaman serisi analizi, doğal dil işleme ve benzeri uygulamalarda tercih edilir [31].

Bu model, iki LSTM katmanı, 2 *Dropout* katmanları, *RepeatVector* katmanı ve bir *Dense*

katmanından oluşmaktadır. Model, *Adam* optimizasyonu ve ortalama karesel hata kaybı ile derlenir. LSTM katmanlarının birinci katmanındaki aktivasyon fonksiyonu *sigmoid*'dir, ve *kernel*, *recurrent* ve *bias* regularizasyonları L2 normu kullanılarak 0.003 oranında uygulanır. Öğrenme oranı 0.005 olarak belirlenmiştir. *Dropout* oranları ise sırasıyla 0.2'dir. *Dropout* katmanı ile rastgele nöronların oluşması engellenip aşırı uyum azaltılmaktadır. *RepeatVector* katmanında, ikinci LSTM katmanının çıkışını, orijinal zaman serisi boyutuna geri döndürmek için kullanılır. Son olarak, *Dense* katmanı, modelin çıkışını oluşturur. Bu katmanın çıkış boyutu, modelin öğrenmeye çalıştığı veri setinin özellik sayısına eşittir.

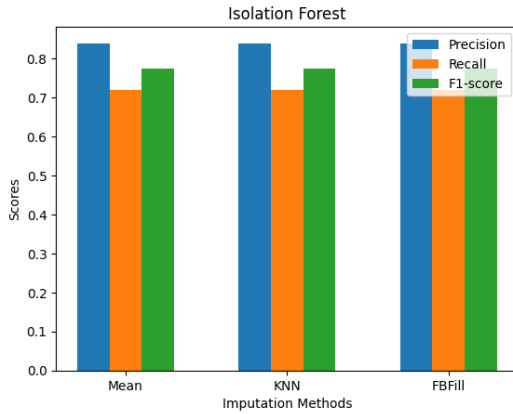
4. BULGULAR VE ANALİZ

Bu başlıkta eğitilen modellerden elde edilen bulgular ve bunların analizi anlatılmıştır. Önce her bir modelin çıktısı verilmiş olup altına da yapılan analiz yazılmıştır.

4.1 Isolation Forest

Tablo 4. IF Sonuçları

	Etiket	P	R	F ₁	A
Ortalama	0	0.76	0.86	0.81	0.79
	1	0.84	0.72	0.78	
K-NN	0	0.76	0.86	0.81	0.79
	1	0.84	0.72	0.78	
FB-Fill	0	0.76	0.86	0.81	0.79
	1	0.84	0.72	0.78	



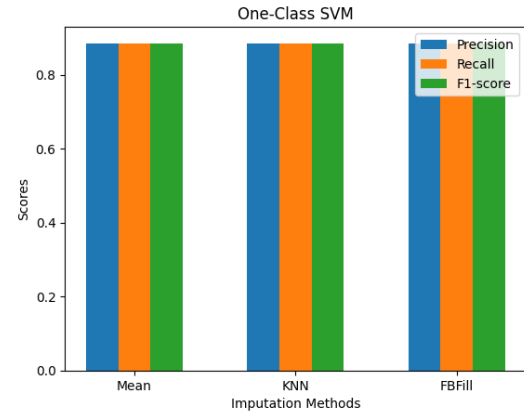
Şekil 9. IF Sonuçları

Tablo 4 ve Şekil 9'dan yola çıkarak bu modelin imputasyon yöntemlerine göre fark göstermeden aynı performansı gösterdiği çıkarılabilir. Model %84 P elde ederek, etiketlediği anomali sınıfındaki örneklerin %84'ünün gerçekten anomali olduğunu doğrulamıştır. Ayrıca, %72 R elde ederek, gerçek anormal örneklerin %72'sini doğru bir şekilde tespit etmiştir. Bu değerlerden yola çıkarak %78 F₁ Skoru hesaplanmıştır.

4.2 One-Class SVM

Tablo 5. 1-C SVM Sonuçları

	Etiket	P	R	F ₁	A
Ortalama	0	0.89	0.89	0.89	0.88
	1	0.89	0.89	0.89	
K-NN	0	0.89	0.89	0.89	0.88
	1	0.89	0.89	0.89	
FB-Fill	0	0.89	0.89	0.89	0.88
	1	0.89	0.89	0.89	



Şekil 10. 1-C SVM Sonuçları

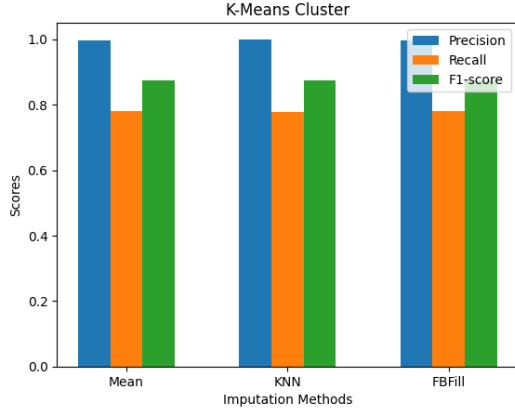
Tablo 5 ve Şekil 10'da görüldüğü üzere bu model anomalileri tespit etmekte hem hassasiyet olarak hemde duyarlılık olarak %89 oranında bir başarı elde etmiştir. P, R ve F₁ oranlarının aynı olması modelin anormal durumları etiketleme yeteneğinin oldukça güçlü olduğunu göstermektedir.

4.3 K-Means Clustering

Tablo 6. K-Means Clustering Sonuçları

	Etiket	P	R	F ₁	A
Ortalama	0	0.82	1	0.90	0.88
	1	1	0.78	0.88	
K-NN	0	0.82	1	0.90	0.88
	1	1	0.78	0.88	

FB-Fill	0	0.82	1	0.90	0.88
	1	1	0.78	0.88	



Şekil 11. K-Means Clustering Sonuçları

Bu modelin sonuçlarından modelin tüm anomalileri tespit edebildiğini, bu anomalilerin %78'inin doğru olduğu çıkartılmaktadır. %88 F₁ skoruyla bu modelinde anormal durumları etiketleme yeteneğinin güçlü olduğu söylenebilir.

4.4 GMM

Diğer modellerin aksine bu modelde imputasyon yöntemine göre sonuçlarda farklılıklar vardır.

Tablo 7. GMM Sonuçları

	Etiket	P	R	F ₁	A
Ortalama	0	0.81	0.97	0.88	0.87
	1	0.96	0.77	0.86	
K-NN	0	0.80	0.96	0.88	0.86
	1	0.95	0.76	0.85	
FB-Fill	0	0.79	0.95	0.90	0.84
	1	0.93	0.75	0.83	

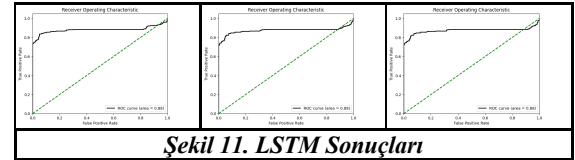
Model ortalama yöntemiyle impute edilmiş veriyle eğitildiğinde %96 hassasiyetle anomali tespit edebilmekte ve bunun %77'sinin doğru tespit edildiği gözükmemektedir. K-NN yöntemiyle işlenmiş verinin beslenmesiyle modelin hassasiyet ve duyarlılığı %1 oranında düşmüştür. FB-Fill yöntemiyle işlenmiş verinin beslenmesi sonucu ise model, anomalilerin %93'ünü tespit edebilmekte ve bunların %75'ini doğru bilmektedir. F₁ skorları karşılaştırıldığında üç imputasyon yönteminin, modelin eğitiminde kayda değer bir değişim yaptığı söylenemez.

4.5 LSTM

Bu model bir derin öğrenme modeli olduğu için eğitimi diğerlerinden farklıdır. Eğitim verisi, aykırı değerlerinden ayıklanıp PCA uygulanmadan modele verilmiştir. Veri seti önce Min-Max yöntemiyle ölçeklendirilmiştir. Daha sonra eğitime başlanmıştır.

Tablo 8. LSTM Sonuçları

	Etiket	P	R	F ₁	ROC
Ortalama	0	0.85	0.95	0.90	0.89
	1	0.94	0.84	0.89	
K-NN	0	0.85	0.93	0.89	0.88
	1	0.92	0.84	0.88	
FB-Fill	0	0.85	0.93	0.89	0.88
	1	0.92	0.84	0.88	



Şekil 11. LSTM Sonuçları

Tablo 8 ve Şekil 11'deki bulgulardan yola çıkarak bu modelin kendi öz niteliklerini çıkarıp anomali tespitinde başarılı olduğu gözlenmektedir. Ortalama imputasyonuyla işlenmiş veriyle beslenen model, %94 hassasiyetle anomali tespiti yapıp bunların %84'ünü doğru bilirken, diğer iki yöntemle işlenmiş veriyle beslenen model ise sadece %2'lik bir kayıpla hassasiyet bakımından geri kalmıştır. F₁ skorları ve ROC alanları göz önüne alındığında bu model anomali tespitinde oldukça başarılıdır.

5. SONUÇLAR VE ÖNERİLER

Bu çalışma kapsamında kullanılan veri seti için eksik değerlerin yerine konulması için uygulanan ortalama, K-NN ve FB-Fill yöntemlerinin, modellerin performanslarına çok büyük bir katkı sağlamadığı çıkarımı yapılabilir. Her bir modelin imputasyon yöntemi gözetmeksizin en iyi F₁ skorları ele alındığında aşağıdaki tablo ortaya çıkmaktadır.

	IF	IC-SVM	K-Means	GMM	LSTM
F₁	0.78	0.89	0.88	0.86	0.89

Bu sonuçlara göre, One-Class SVM ve LSTM modellerinin, veri setindeki anormal durumları etiketleme konusunda diğer modellere göre daha başarılı bir performans sergilediği görülmektedir.

1 Model güncelleme ve ayarlama süreçleri üzerinde
2 çalışmak, mevcut modellerin daha da
3 iyileştirilmesine katkı sağlayabilir. Ayrıca model
4 sonuçlarının iş uygulamalarına entegrasyonunu
5 sağlamak ve gerçek dünya senaryolarında nasıl
6 kullanılacağı düşünülebilir. Bu, anomali tespiti
7 modellerinin etkili bir şekilde implemente edilerek,
8 iş süreçlerinde olası anormal durumları önceden
9 belirleme yeteneğini artırabilir.

10 REFERANSLAR

11
12 [1] A. B. Nassif, M. A. Talib, Q. Nasir and F. M.
13 Dakalbab, "Machine learning for anomaly detection:
14 A systematic review", IEEE Access, vol. 9, pp.
15 78658-78700, 2021.
16 [2] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly
17 detection: A survey," ACM Comput. Surv., vol. 41,
18 no. 3, pp. 71–97, 2009,
19 doi:10.1145/1541880.1541882.
20 [3] M. Injadat, F. Salo, A. B. Nassif, A. Essex, and A.
21 Shami, "Bayesian optimization with machine
22 learning algorithms towards anomaly detection," in
23 Proc. IEEE Global Commun. Conf. (GLOBECOM),
24 Dec. 2018, pp. 1–6, doi:
25 10.1109/GLOCOM.2018.8647714.
26 [4] Anomaly Detection in Helicopter,
27 [https://www.kaggle.com/datasets/nagasai524/anoma-](https://www.kaggle.com/datasets/nagasai524/anomaly-detection-in-helicopter/data)
28 [ly-detection-in-helicopter/data](https://www.kaggle.com/datasets/nagasai524/anomaly-detection-in-helicopter/data).
29 [5] H. Izakian and W. Pedrycz. 2013. Anomaly detection
30 in time series data using a fuzzy c-means clustering.
31 In Proceedings of the 2013 Joint IFSA World
32 Congress and NAFIPS Annual Meeting. IEEE,
33 Edmonton, Alberta, Canada, 1513–1518.
34 [6] R. J. Hyndman, E. Wang, and N. Laptev. 2015.
35 Large-Scale Unusual Time Series Detection. In
36 Proceedings of the 15th IEEE International
37 Conference on Data Mining Workshop (ICDMW
38 '15). IEEE, Atlantic City, NJ, USA, 1616–1619.
39 [7] L. Xiong, H.-D. Ma, H.-Z. Fang, K.-X. Zou and D.-
40 W. Yi, "Anomaly detection of spacecraft based on
41 least squares support vector machine", Proc.
42 Prognostics Syst. Health Management Conf., pp. 1-6,
43 May 2011..
44 [8] H. N. Akouemo and R. J. Povinelli. 2016.
45 Probabilistic anomaly detection in natural gas time
46 series data. Int. J. Forecast. 32, 3 (2016), 948–956.
47 [9] N. Laptev, S. Amizadeh, and I. Flint. 2015. Generic
48 and Scalable Framework for Automated Time-series
49 Anomaly Detection. In Proceedings of the 21th ACM
50 SIGKDD International Conference on Knowledge
51 Discovery and Data Mining (KDD '15). ACM,
52 Sydney, NSW, Australia, 1939–1947.
53 [10] A. Nanduri and L. Sherry, "Anomaly Detection in
54 Aircraft Data Using Recurrent Neural Networks
55 (RNN)", Proceedings of Integrated Communications
56 Navigation and Surveillance, pp. 5C2-1-5C2-8, April
57 2016.
58 [11] Yan, S., Shao, H., Xiao, Y., Liu, B., & Wan, J.
59 (2023). Hybrid robust convolutional autoencoder for

60 unsupervised anomaly detection of machine tools
61 under noises. Robotics and Computer-Integrated
62 Manufacturing, 79, Article 102441.
63 [12] Vos K, Peng Z, Jenkins C, Shahriar MR, Borghesani
64 P, Wang W. Vibration-based anomaly detection
65 using LSTM/SVM approaches. Mech Syst Signal
66 Process 2022; 169:108752.
67 <https://doi.org/10.1016/j.ymssp.2021.108752>.
68 [13] Liu, C., & Gryllias, K. . (2021). A Deep Support
69 Vector Data Description Method for Anomaly
70 Detection in Helicopters. PHM Society European
71 Conference, 6(1), 9.
72 <https://doi.org/10.36001/phme.2021.v6i1.2957>.
73 [14] Ammann, Oliver et al. "Anomaly Detection And
74 Classification In Time Series With Kervolutional
75 Neural Networks." ArXiv abs/2005.07078 (2020).
76 [15] M. Jones, D. Nikovski, M. Imamura, and T. Hirata.
77 2016. Exemplar learning for extremely efficient
78 anomaly detection in real-valued time series. Data
79 Min.Knowl. Discov. 30, 6 (2016), 1427–1454.
80 [16] Chandola, Varun; Cheboli, Deepthi; Kumar, Vipin.
81 (2009). Detecting Anomalies in a Time Series
82 Database. Retrieved from the University of
83 Minnesota Digital Conservancy,
84 <https://hdl.handle.net/11299/215791>.
85 [17] Mohamed Ben Slimene, Mohamed-Salah
86 Ouali, Anomaly Detection Method of Aircraft System
87 using Multivariate Time Series Clustering and
88 Classification Techniques, IFAC-
89 PapersOnLine, Volume 55, Issue 10, 2022, Pages
90 1582-1587, ISSN 2405-
91 8963, <https://doi.org/10.1016/j.ifacol.2022.09.616>.
92 [18] Iliopoulos, Anastasios, et al. "Detection of Anomalies
93 in Multivariate Time Series Using Ensemble
94 Techniques." 2023 IEEE Ninth International
95 Conference on Big Data Computing Service and
96 Applications (BigDataService). IEEE, 2023.
97 [19] G. Li, A. Rai, H. Lee, A. Chattopadhyay,
98 "Operational anomaly detection in flight data using a
99 multivariate gaussian mixture model". In proceedings
100 of the 10th Annual Conference of the Prognostics and
101 Health Management Society, 2018. doi:
102 10.36001/phmconf.2018.v10i1.474.
103 [20] Z. Ghrib, R. Jaziri and R. Romdhane, "Hybrid
104 approach for Anomaly Detection in Time Series
105 Data," 2020 International Joint Conference on Neural
106 Networks (IJCNN), Glasgow, UK, 2020, pp. 1-7, doi:
107 10.1109/IJCNN48605.2020.9207013.
108 [21] A. Gaev and A. Lantsberg, "A Method to Handle
109 Unstable Time Series in Anomaly Detection
110 Problem," 2021 3rd International Conference on
111 Control Systems, Mathematical Modeling,
112 Automation and Energy Efficiency (SUMMA),
113 Lipetsk, Russian Federation, 2021, pp. 548-552, doi:
114 10.1109/SUMMA53307.2021.9632243.
115 [22] J. Shi, G. He and X. Liu, "Anomaly Detection for Key
116 Performance Indicators Through Machine Learning,"
117 2018 International Conference on Network
118 Infrastructure and Digital Content (IC-NIDC),
119 Guiyang, China, 2018, pp. 1-5, doi:
120 10.1109/ICNIDC.2018.8525714.

- [23] H. C. Mandhare and S. R. Idate, "A comparative study of cluster based outlier detection, distance based outlier detection and density based outlier detection techniques," 2017 International Conference on Intelligent Computing and Control Systems (ICICCS), Madurai, India, 2017, pp. 931-935, doi: 10.1109/ICCONS.2017.8250601.
 - [24] Y. Jiang, N. Le, Y. Zhang, Y. Zheng and Y. Jiao, "Research on the Flight Anomaly Detection During Take-off Phase Based on FOQA Data," 2019 CAA Symposium on Fault Detection, Supervision and Safety for Technical Processes (SAFEPROCESS), Xiamen, China, 2019, pp. 756-760, doi: 10.1109/SAFEPROCESS45799.2019.9213422.
 - [25] Bishop, Christopher (2006). Pattern recognition and machine learning. Berlin: Springer. ISBN 0-387-31073-8.
 - [26] Liu, Fei Tony. "First Isolation Forest implementation on Sourceforge".
 - [27] Liu, Fei Tony; Ting, Kai Ming; Zhou, Zhi-Hua (December 2008). "Isolation Forest". 2008 Eighth IEEE International Conference on Data Mining. pp. 413–422. doi:10.1109/ICDM.2008.17. ISBN 978-0-7695-3502-9. S2CID 6505449.
 - [28] Liu, Fei Tony; Ting, Kai Ming; Zhou, Zhi-Hua (December 2008). "Isolation-Based Anomaly Detection". ACM Transactions on Knowledge Discovery from Data. 6: 3:1–3:39. doi:10.1145/2133360.2133363. S2CID 207193045.
 - [29] Hastie, T.; Tibshirani, R.; Friedman, J. (2009). "Unsupervised Learning". The Elements of Statistical Learning (2nd ed.). Springer.
 - [30] Bishop, C. M. (2006). Pattern Recognition and Machine Learning. Springer.
 - [31] Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. Neural computation, 9(8), 1735-1780.
- github:
https://github.com/frkanyilmaz2/machine_learning/blob/main/main.ipynb