

TOBB EKONOMİ VE TEKNOLOJİ ÜNİVERSİTESİ
FEN BİLİMLERİ ENSTİTÜSÜ

YOĞUN BAKIM HASTALARININ MORTALİTE ve HASTANEDE KALMA
SÜRELERİNİN DERİN ÖĞRENME YÖNTEMLERİ ile TAHMİNİ



DOKTORA TEZİ

Batuhan BARDAK

Bilgisayar Mühendisliği Anabilim Dalı

Tez Danışmanı: Doç. Dr. Mehmet TAN

TEMMUZ 2022

TEZ BİLDİRİMİ

Tez içindeki bütün bilgilerin etik davranış, ve akademik kurallar çerçevesinde elde edilerek sunulduğunu, alıntı yapılan kaynaklara eksiksiz atıf yapıldığını, referansların tam olarak belirtildiğini ve ayrıca bu tezin TOBB ETÜ Fen Bilimleri Enstitüsü tez yazım kurallarına uygun olarak hazırlandığını bildiririm.

Batuhan BARDAK

ÖZET

Doktora Tezi

YOĞUN BAKIM HASTALARININ MORTALİTE ve HASTANEDE KALMA SÜRELERİNİN DERİN ÖĞRENME YÖNTEMLERİ ile TAHMİNİ

Batuhan BARDAK

TOBB Ekonomi ve Teknoloji Üniversitesi
Fen Bilimleri Enstitüsü
Bilgisayar Mühendisliği Anabilim Dalı

Tez Danışmanı: Doç. Dr. Mehmet TAN

Tarih: TEMMUZ 2022

Günümüzde dijital dönüşüm hızının artmasıyla beraber fiziksel olarak saklanan verilerin elektronik ortamlara aktarılmasında hız kazanmıştır. Bu durum, birçok alana olduğu gibi sağlık alanına da doğrudan yansımıştır. Geçmişte fiziksel olarak saklanan hastaya ait kayıtlar bu sayede dijital ortamlara geçirilmiştir. Dijital ortama aktarılan hastaya ait demografik bilgiler, laboratuvar sonuçları, yaşamsal gözlem verileri, klinik notlar, tanı kodları ve benzeri birçok veri Elektronik Sağlık Kaydı (ESK) olarak tanımlanmaktadır. Sağlık alanındaki dijital dönüşüme ek olarak, derin öğrenme yöntemlerine olan geniş ilgi, araştırmacıları, finans, sosyal medya, siber güvenlik gibi birçok alanda yapay zeka yöntemlerini kullanmaya teşvik etmektedir. Elektronik sağlık kayıtlarının araştırmacılar için kullanılabilir hale gelmesiyle birlikte, bu veri setlerini kullanarak derin öğrenme modelleri geliştirmeye olan ilgi artmaktadır. Tez kapsamında yapılan deneylerde, günümüzdeki en popüler ve erişilebilir elektronik sağlık kayıt veri seti olan Medical Information Mart for Intensive Care (MIMIC-III) kullanılmıştır. Yoğun bakımda yatan hastaların, yaşamsal gözlem verilerini ve diğer klinik bilgilerini ölçerek, hastaların mevcut sağlık durumlarını anlamlandırmak ve gelecek sağlık durumlarını tahmin etmek önemli bir problemdir. Tez kapsamında, hastaların hastane içinde ve yoğun bakımda mortalite ihtimalleri ile yoğun bakımda 3 ve 7 günden fazla kalıp kalmayacakları çok-kipli derin öğrenme tabanlı yöntemler ile tahmin edilmiştir.

Gerçekleştirilen çalışma üç ana bölüme ayrılmıştır. İlk bölümde, yoğun bakımda yatan hastalara ait yaşamsal gözlem verileri, laboratuvar sonuçları gibi özniteliklere ek olarak hastalara ait klinik notlar da model eğitime dahil edilmiş ve modelin

klinik problemleri tahmin etme başarısı artırılmaya çalışılmıştır. İkinci bölümde, klinik notların doğrudan kullanılması yerine, varlık isim tanıma yöntemi ile notlar içerisinden medikal terimlerin çıkartılması sağlanmıştır. Elde edilen medikal terimlerin, mortalite ve yoğun bakımda kalma süresi tahmini problemlerine etkisi araştırılmıştır. Yapılan son çalışmada ise, hastaların zaman serisi özniteliklerine ilave olarak, hastaların yoğun bakımda kaldıkları süre boyunca kullandıkları ilaçların moleküler temsilleri kullanılmış, ve klinik problemlerin tahminine etkisi üzerine deneyler yapılmıştır. Ek olarak, bu çalışma sonunda, hastanede mortalite tahmini için eğitilen modelin açıklanabilirliğini arttırmak amacıyla SHapley Additive exPlanations (SHAP) yöntemi kullanılmıştır. SHAP yönteminin çıktısı, zaman-serisi ve klinik ilaç özniteliklerinin model üzerindeki etkisininin daha derin bir analizinin yapılmasına olanak sağlamaktadır.

MIMIC-III veri seti içerisinde hastaya ait farklı veri türlerinin bir arada bulunması, tez kapsamında yapılan deneylerde bu veri türlerinin bir arada kullanılabilmesine ve farklı deneylerin gerçekleştirilebilmesine olanak sağlamıştır. Farklı veri türlerini aynı model içerisinde kullanabilmek için çok-kipli derin öğrenme tabanlı yöntemler önerilmiştir. Yapılan deney sonuçları incelendiğinde, zaman-serisi özniteliklerin yanısıra hastaya ait klinik notların, medikal terimlerin ve ilaç bilgilerinin modele girdi olarak verilmesinin, klinik problemlerin başarımına olumlu yönde etki ettiği görülmüştür.

Anahtar Kelimeler: Elektronik sağlık kaydı, Derin öğrenme, Çok-kipli modeller, Doğal dil işleme, Açıklanabilir yapay zeka

ABSTRACT

Doctor of Philosophy

PREDICTION OF MORTALITY AND LENGTH OF STAY OF ICU PATIENTS WITH DEEP LEARNING

Batuhan BARDAK

TOBB University of Economics and Technology
Institute of Natural and Applied Sciences
Department of Computer Engineering

Supervisor: Assoc. Prof. Mehmet TAN

Date: July 2022

The increase in the speed of digital transformation has accelerated the transfer of physically stored data to electronic media data. Healthcare is one of many areas that is impacted by these transformation. Electronic health records (EHR) is the general term for the data that is associated with a patient's whole health journey including demographic information, laboratory test results, vital signs, clinical notes, diagnosis codes, and related data. In addition to the digital transformation in healthcare, the widespread interest in machine/deep learning encourages the researchers to apply artificial intelligence to several different domains such as finance, social media, and cyber security. With the EHR data becoming available for researchers, there has been an increasing interest in using it with deep learning algorithms. Within the scope of this study, we use the most popular and publicly available EHR dataset, Medical Information Mart for Intensive Care (MIMIC-III). Understanding the health condition of the patient by observing the clinical measurements, and laboratory tests, and predicting the condition of patients during their intensive care unit (ICU) stay is a vital problem. In this study, two different common risk prediction tasks, mortality (in-hospital & in-ICU), and length of ICU stay ($LOS > 3$, $LOS > 7$) are researched.

The interest of this work is divided into three parts. In the first part, we use the clinical notes besides the time-series features such as vital signs and laboratory test results to

improve the model predictions. In the second part, instead of using clinical notes directly, we extract medical entities from clinical notes by clinical named entity recognition (NER) model and use them as additional features besides time-series features to improve proposed model predictions. In the last study, we argue the integration of structured time-series data and molecular representations of the drugs which are prescribed to patients in ICU. Several experiments are conducted to investigate the effect of clinical drugs on mortality and LOS problem predictions. Additionally, the SHapley Additive exPlanations (SHAP) is applied to increase the interpretability of the in-hospital mortality model and to investigate the relationship between the mortality and the time-series and clinical drug features. The output of the SHAP method allows us to make a deeper analysis of the effect of time-series and clinical drug features.

Since MIMIC-III contains rich information with multiple modalities of data, we apply a multimodal learning approach to handle the heterogeneous nature of the data. The experimental results indicate a promising increase in performance on clinical tasks when the clinical notes, medical entities or clinical drug informations are used with time series features in a multimodal approach.

Keywords: Electronic health record, Deep learning, Multimodal learning, Natural language processing, Explainable artificial intelligence

TEŞEKKÜR

Doktora öğrenim hayatı ve tez çalışmalarım boyunca değerli bilgi birikimi ve tecrübesiyle bana yol gösteren, zamanını, desteğini ve hoşgörüsünü benden esirgemeyen kıymetli hocam Doç. Dr. Mehmet TAN'a teşekkürü bir borç bilirim. Ayrıca, değerli vakitlerini ayırarak yaptıkları yorumlarla tezime katkı sağlayan tez izleme komitesi üyeleri Prof. Dr. Pınar KARAGÖZ ve Doç. Dr. Murat ÖZBAYOĞLU'na, tez savunmam sırasında yaptıkları değerli yorumlar ile tezimi iyileştirmeme yardımcı olan Prof. Dr. Ferda Nur ALPASLAN ve Prof. Dr. Osman ABUL'a, öğrenim hayatım boyunca tecrübelerinden faydalandığım TOBB Ekonomi ve Teknoloji Üniversitesi Bilgisayar Mühendisliği Bölümü'nün değerli tüm öğretim üyelerine teşekkürlerimi sunarım.

Doktora öğrenimim süresince araştırma bursu vererek bana destek olan TOBB Ekonomi ve Teknoloji Üniversitesi'ne ve 120E173 nolu TÜBİTAK 1002 projesi ile bu tez çalışmamı desteklediği için TÜBİTAK'a ayrıca teşekkür ederim.

Doktora sürecimde değerli fikirleri ve yardımlarıyla bana destek olan sevgili dostlarım Mehmet Saygın SEYFİOĞLU ve Abdullah Serdar YONAR'a, bu zorlu süreçteki destekleri için STM Büyük Veri ekibindeki çalışma arkadaşlarıma ve yöneticilerime, ve son olarak hayatı eğlenceli kılan diğer tüm arkadaşlarıma çok teşekkür ederim.

Sadece bu tez çalışması boyunca değil, tanıştığımız günden beri bana olan sevgisini ve desteğini hiç eksik etmeyen, elinden gelen her türlü yardımını sunan nişanlım Sıla ŞİBİL'e çok teşekkür ederim.

Son olarak, beni büyüten, bana güç veren ve bu günlere gelmemde en büyük pay sahibi olan değerli annem Nilgün BARDAK, babam Metin BARDAK ve her zaman yanımda olan kardeşim Benan BARDAK'a çok teşekkür ederim.

İÇİNDEKİLER

	<u>Sayfa</u>
ÖZET	iv
ABSTRACT	vi
TEŞEKKÜR	viii
İÇİNDEKİLER	ix
ŞEKİL LİSTESİ	xi
ÇİZELGE LİSTESİ	xiii
KISALTMALAR	xv
1. GİRİŞ	1
1.1 Problemin Tanımı.	3
1.2 Arastırma Motivasyonu ve Önerilen Çözüm Yöntemi	3
1.3 Arastırma Zorlukları.	4
1.4 Tezin Planı.	5
1.5 Tezin Katkıları.	5
1.6 Tezin Düzeni	7
2. İLGİLİ ÇALIŞMALAR	9
3. VERİ KÜMESİ	15
3.1 MIMIC-III Veri Kümesi	15
3.1.1 MIMIC-III Veri setine erişim	16
3.1.2 MIMIC-III Tablo detayları	17
3.2 MIMIC-III Ön İşleme Adımları	19
4. KULLANILAN YÖNTEMLER.....	27
4.1 ESK Verilerinin Analizinde Derin Öğrenme Yöntemleri	27
4.1.1 Tekrarlamalı yapay sinir ağı(RNN)	27
4.1.2 Uzun kısa vadeli hafıza ağları(LSTM)	28
4.1.3 Geçitli tekrarlayan sinir ağları(GRU)	31
4.1.4 Evrimsel sinir ağları(CNN)	33
4.2 Kelime Temsil Yöntemleri.....	36
4.2.1 Word2Vec.....	36
4.2.2 FastText.....	39
4.2.3 Doc2Vec.....	39
4.2.4 Bidirectional Encoder Representations from Transformers (BERT) .	40
4.2.5 Clinical BERT.....	43
4.2.6 Sentence-BERT	44
4.3 Klinik Varlık İsimlerinin Tanımlanması.	45
4.4 İlaçların Temsili	46
4.4.1 Extended-Connectivity Fingerprints (ECFP)	47

4.4.2 Molecular Access System (MACCS)	48
4.4.3 Mol2Vec	48
4.4.4 Smiles-Transformer	49
4.5 SHapley Additive exPlanations (SHAP)	50
4.6 Performans Metrikleri	53
4.7 Kullanılan Kütüphaneler ve Çalışma Ortamı	55
5. ZAMAN SERİSİ ve MEDİKAL TERİMLER ile TAHMİN ETME.....	57
5.1 Motivasyon	57
5.2 Önerilen Yöntem.....	58
5.3 Deneysel Sonuçlar.....	64
5.4 Değerlendirme	66
6. ZAMAN SERİSİ ve KLİNİK NOTLAR ile TAHMİN ETME.....	69
6.1 Motivasyon	69
6.2 Önerilen Yöntem.....	69
6.3 Deneysel Sonuçlar.....	72
6.4 Değerlendirme	74
7. ZAMAN SERİSİ ve İLAÇ TEMSİLLERİ ile TAHMİN ETME.....	77
7.1 Motivasyon	77
7.2 Önerilen Yöntem.....	78
7.3 Deneysel Sonuçlar.....	83
7.4 Değerlendirme	85
8. SONUÇ ve ÖNERİLER.....	93
8.1 Gelecek Çalışmalar için Öneriler.....	95
KAYNAKLAR.....	96
ÖZGEÇMİŞ	109

ŞEKİL LİSTESİ

	<u>Sayfa</u>
Şekil 1.1: Tezin temel adımları.	6
Şekil 3.1: Çeşitli sağlık verilerinin oluşum şekilleri.	16
Şekil 3.2: MIMIC-III genel mimarisi.	17
Şekil 3.3: Örnek hasta ziyaret ve kayıt defteri.	19
Şekil 3.4: MIMIC-Extract temel veri ön işleme adımları	24
Şekil 3.5: MIMIC-Extract veri ön işleme sonrası, hastalarının yaş dağılımı	24
Şekil 3.6: MIMIC-Extract veri ön işleme sonrası, hastalarının yoğun bakımda kalma sürelerinin dağılımı	25
Şekil 4.1: Yinelemeli sinir ağı mimarisi.	29
Şekil 4.2: Uzun kısa vadeli hafıza ağı mimarisi.	30
Şekil 4.3: LSTM mimarisi hücre içi yapısı	31
Şekil 4.4: GRU hücre içi yapısı	33
Şekil 4.5: CNN’de evrişim işlemi.	35
Şekil 4.6: CNN’de dolgulama (padding) işlemi.	36
Şekil 4.7: CNN’de örnekleme (pooling) işlemi.	37
Şekil 4.8: Metinsel veri üzerinde 1-boyutlu evrişim ve maksimum örnekleme işlemi.	38
Şekil 4.9: Skip-Gram yönteminin gösterimi.	38
Şekil 4.10: CBOW yönteminin gösterimi.	39
Şekil 4.11: Klinik alana özgü Word2Vec eğitimi.	40
Şekil 4.12: Doc2Vec yöntemi içerisindeki farklı öğrenme mimarileri olan PV-DM ve PV-DBOW metotlarının gösterimi.	41
Şekil 4.13: Transformer, $BERT_{BASE}$ ve $BERT_{LARGE}$ yöntemlerinin mimari büyüklük karşılaştırılması.....	42
Şekil 4.14: BERT modelinin ön eğitim aşamasının gösterimi.	43
Şekil 4.15: BERT modeli ile cümle benzerlik görevi mimarisi.	45
Şekil 4.16: Sentence-BERT mimarisi.	46
Şekil 4.17: med7 çalışması kapsamında eğitilen klinik VİT modelinin, MIMIC-III içerisindeki bir klinik not cümlesi üzerindeki örnek çıktısı.....	47
Şekil 4.18: ECFP yöntemi örnek gösterimi.	48
Şekil 4.19: Mol2Vec yöntemi örnek gösterimi.	49
Şekil 4.20: SMILES-Transformer yöntemi ön-eğitim modeli.	50
Şekil 4.21: SHAP öznitelik kombinasyon gösterimi.	51
Şekil 4.22: Karışıklık matrisi.	53
Şekil 4.23: Örnek Kararlılık-Duyarlılık (PR) grafiği.	55
Şekil 4.24: Örnek ROC grafiği.	55

Şekil 5.1: Medikal varlık isim vektörlerini öğrenmek için önerilen yöntem. (1) Klinik notlar içerisinden çıkartılan medikal varlıklar sürekli değerler içeren vektörlere dönüştürürler. Daha sonra ise öğrenilmiş temsillerin ortalaması alınmıştır. (2) Eğer klinik notlar içerisindeki kelimelerden biri medikal varlık tiplerinden herhangi birine girmiyor ise o kelimeler klinik notlar içerisinden silinmiştir. Ardından, veri ön işleme uygulanmış, klinik notlar üzerinden Doc2Vec yöntemi eğitilmiştir	62
Şekil 5.2: Hastane içi mortalite, YBÜ mortalite, Hastanede Kalma Süresi >3, ve Hastanede Kalma Süresi >7 klinik problemleri için önerilen yöntem özetleri. MIMIC-III içerisinde zamana bağlı özniteliklerin hazırlanması için MIMIC-Extract çalışması kullanılmış, ve öznitelikler GRU algoritmasına girdi olarak verilmiştir. Aynı zamanda klinik notlara veri ön işleme adımları uygulanmış, ve medikal terimlerin çıkartılması için med7 yöntemi kullanılmıştır. Medikal terim temsilleri içerisinden öznitelik çıkarımı için 1D CNN yöntemi önerilmiştir. Son olarak ise 4 farklı sınıflandırma problemini tahminleyebilmek için çıkartılan öznitelikler birleştirilmiş, ve 3 katmanlı yapay sinir ağına girdi olarak verilmiştir	63
Şekil 6.1: Mortalite ve Yoğun Bakımda Kalma Süresi temelli dört ayrı klinik problemin tahmini için önerilen model mimarisi.	73
Şekil 7.1: MIMIC-III içerisinden klinik ilaç isimlerini çıkartan yöntem. Çıkartılan ve ön işlemden geçirilen klinik ilaçlar PubChem içerisinde bulunduktan sonra farklı moleküller temsillere dönüştürülmektedir.....	81
Şekil 7.2: Önerilen çok-kipli model mimarisinin özeti. İlk olarak, zaman-serisi ve ilaç verileri ön işlemden geçirilmiştir. İkinci olarak, öznitelik çıkarımı için, zaman-serisi özniteliklere GRU modeli, klinik ilaçlara ise 1D CNN modeli uygulanmıştır. Daha sonra ise, çıkartılan öznitelikler birleştirilerek yapay sinir ağına verilerek, 4 farklı klinik problem için tahmin yapılmıştır. Son olarak ise, SHAP yöntemi kullanılarak, hastane içi mortalite problemi ile önemli zaman-serisi öznitelikler ve ilaçlar arasında ilişki yakalanmaya çalışılmıştır.....	83
Şekil 7.3: SHAP uygulamasının çalışma içerisindeki kullanımı	84
Şekil 7.4: SHAP çıktılarını önemli özniteliklerin seçimi. Deneysel sonuçlar üç farklı figürde gösterilmiştir. Figürlerdeki X-ekseni (x,y) formatında olmak üzere çıkartılan öznitelik sayılarını temsil etmektedir. x sembolü çıkartılan zaman-serisi öznitelik sayısını temsil ederken, y ise çıkartılan klinik ilaç sayısını temsil etmektedir.....	89

ÇİZELGE LİSTESİ

	<u>Sayfa</u>
Çizelge 3.1: MIMIC-III tabloları ve kısa açıklamaları	20
Çizelge 3.2: MIMIC-Extract çalışması ile çıkartılan ve deneylerde kullanılan 104 adet zaman-serisi öznitelik.	23
Çizelge 3.3: MIMIC-III veri setinin, MIMIC-Extract uygulanmadan önceki ve sonraki veri istatistikleri.	24
Çizelge 5.1: MIMIC-III veri setinin ve bu çalışmada kullanılan veri setinin istatistikleri.	60
Çizelge 5.2: İlk sütun medikal terimlerin tiplerini belirtirken, ikinci sütun derlem içerisindeki klinik notlarda bu terimlerin isimlerinin kaçar kez geçtiği göstermektedir. Üçüncü sütun, bu ifadelerin tekil olarak kaçar kez gösterirken, son sütunda ise her bir terim için örnek ifade verilmiştir.	61
Çizelge 5.3: Temel yöntemlerin performans karşılaştırılması. 4 ayrı klinik problemin her biri için, AUC, AUPRC ve F1 skorlarının ortalaması ve standart sapma değerleri raporlanmıştır	65
Çizelge 5.4: Önerilen yöntem ile diğer en iyi temel yöntemlerin karşılaştırılması.	67
Çizelge 6.1: Önerilen model performansının, temel olarak alınan modelle karşılaştırılması. Her problem ve metrik için en yüksek skorlar vurgulanmıştır.	71
Çizelge 6.2: MIMIC-III ve bu çalışmada kullanılan veri seti istatistikleri.	72
Çizelge 6.3: Çeşitli medikal kısaltmalar ve bu kısaltmaların dönüştürülen halleri.	73
Çizelge 6.4: Önerilen model performansının, temel olarak alınan modelle karşılaştırılması. Her problem ve metrik için en yüksek skorlar vurgulanmıştır.	74
Çizelge 7.1: MIMIC-III ve bu çalışmada kullanılan veri setinin istatistikleri.	79
Çizelge 7.2: Reçeteli ilaçlar (Prescription) tablosundaki örnek ilaç ismi, genel ismi ve NDC (National Drug Code) örneği.	79
Çizelge 7.3: Örnek klinik ilaç isimleri ve veri ön işleme adımlarından sonraki versiyonu.	80
Çizelge 7.4: Klinik problemler için etiket dağılımları	82
Çizelge 7.5: Dört klinik problem için alınan basarımların ortalama sonuçları ve birbirleri ile karşılaştırılması	86
Çizelge 7.6: Model çıktısına en çok ve en az katkıyı veren ilaç isimlerinin listesi.	89
Çizelge 7.7: Model çıktısına en çok ve en az katkıyı veren zaman-serisi öznitelik listesi.	90

KISALTMALAR

AUPRC: Area Under Precision-Recall Curve
AUPRC: Area Under the Receiver Operating Characteristic Curve
BERT: Bidirectional Encoder Representations from Transformers
BPTT: Backpropagation Through Time
CBOW: Continuous Bag of Words
CITI: Collaborative Institutional Training Initiative
CNN: Convolutional Neural Network
ECFP: Extended-Connectivity Fingerprints
ESK: Elektronik Sağlık Kaydı
FN: False Negative
FP: False Positive

GKS: Glasgow Koma Skalası
GRU: Gated Recurrent Unit
HIPAA: Health Insurance Portability and Accountability Act
ICD: International Classification of Disease
LSTM: Long Short-Term Memory
MACCS: Molecular Access System
MIMIC: Medical Information Mart for Intensive Care
MLM: Masked Language Modeling
n2c2: National NLP Clinical Challenge
NDC: National Drug Code
NER: Named Entity Recognition
NLP: Natural Language Processing
NSP: Next Sentence Prediction
ReLU: Rectified Linear Unit
RNN: Recurrent Neural Network
ROC: Receiver Operating Characteristics
SBERT: Sentence BERT
SHAP: SHapley Additive exPlanations
SMILES: Simplified Molecular Input Line Entry System
TN: True Negative
TP: True Positive
VİT: Varlık İsim Tanıma
XAI: Explainable Artificial Intelligence
YBÜ: Yoğun Bakım Ünitesi

1. GİRİŞ

Hastaların, hastane veya sağlık kuruluşu ziyaretlerinde toplanan verilerine sağlık kaydı ismi verilmektedir. Tarihsel olarak incelendiğinde, eski zamanlardan beri doktorlar, hastalar için gözlemlerini kağıda aktararak sağlık kayıtlarını oluşturmuşlardır. Zamanla artan dünya nüfusu ile, kağıtlarda/dosyalarda tutulan bu sağlık kayıtlarını yönetilemez hale gelmiştir. Teknolojinin gelişmesi ve dijitalleşmenin yaygınlaşması ile dosyalarda tutulan sağlık kayıtlarından Elektronik Sağlık Kaydı(ESK) sistemlerine geçiş başlamıştır. Geleneksel olarak fiziksel bir şekilde saklanan sağlık kayıtlarından ESK'lara geçişte, bu verinin dijital ortamda nasıl tutulacağı ve kullanılacağına dair detaylı araştırmalar ve geliştirmeler yapılmaktadır. Dijitalleşen sağlık kayıtlarının saklanması, paylaşılması, diğer işlemlerinin regüle edilebilmesi ve yasalarının tasarımı için ise 1996 yılında Amerika Birleşik Devletleri'nde HIPAA (Health Insurance Portability and Accountability Act) kurumu kurulmuştur [1]. HIPAA, kişisel sağlık verilerinin kullanım kuralları, bu verilere kimler tarafından erişilebileceği, nasıl paylaşılacağı ve gizliliği (Protected Health Information, PHI) gibi konular üzerinde çalışmakta ve gerekli regülasyonları getirmektedir. Aynı zamanda sağlık verisinin dijital ortamda nasıl saklanacağına yönelik regülasyonlar koyarak, bu alanda çeşitli standartlar yaratmaya çalışmaktadır. Bu gelişmeler ile birlikte hastaneler içerisinde yaygın bir şekilde kullanılmaya başlayan ESK'lar, hastanın geçmiş klinik bilgilerini tutan dijital veri kayıtlarıdır. ESK'lar sayesinde, hastaya ait geçmiş sağlık verileri hem toplu bir biçimde saklanabilir hem de farklı sağlık kuruluşları arasında erişilebilir hale gelmektedir. Sağlık çalışanlarının bu bilgilere erişerek hastanın mevcut durumu hakkında daha iyi analizler yapabilmesine de olanak sağlayan ESK'lar, aynı zamanda yapay öğrenme ve derin öğrenme teknikleri ile birçok problem üzerinde çalışma yapılmasına da olanak sağlamaktadır. ESK'lar ile beraber oluşan büyük heterojen sağlık verisi içerisinde, hastaya ait demografik bilgiler, tedavi bilgileri, klinik notlar ve raporlar, yaşamsal gözlem verileri, medikal kodlar ve medikal görüntüler bulunmaktadır. Büyük ve çeşitli verilerin bir arada bulunması ile sağlık alanında farklı klinik problemlerin yapay zeka temelli çözümleri aranmaya başlanmıştır [2, 3, 4, 5]. Sağlık alanında yapay zeka uygulamalarının kullanılmasının temel iki amacı, sağlık alanındaki maliyetleri azaltmaya çalışmak ve hastalara daha hızlı ve kaliteli sağlık hizmeti sunmaktır. Bunun temel sebebi, sağlık sektörünün hem devletlerin

hem de özel sektörün en çok harcama yaptığı alanların başında gelmesidir. Amerika Birleşik Devletleri'nde yıllık 30 milyara yakın sağlık talebi olmakla beraber, Centers for Medicare & Medicaid Services (CMS)'in raporuna göre 2019 yılındaki sağlık giderleri yaklaşık olarak 3,24 trilyon dolar olmuştur [6]. Bu gider maliyeti, 2022 yılında teknoloji alanında dünyanın en büyük 4 firması olarak kabul edilen Apple, Google, Microsoft ve Amazon'un piyasa değerleri toplamına eşdeğer olduğu görülmektedir. Bir başka önemli husus ise sağlık alanına harcanan paranın büyüklüğü ile beraber israf edilen paranın büyüklüğüdür. Sharank v.d. [6] tarafından yapılan çalışmada yaklaşık 935 milyar dolarlık bütçenin israf edildiği üzerine bulgular okuyucu ile paylaşılmıştır. İsrاف olarak tespit edilen konular ise doktor/hastane hataları, tedavi planlama hataları, gereksiz/yanlış, tedavi uygulamaları, kaçakçılık ve dolandırıcılık (sahte fatura vb.), idari yönetsel karmaşıklıklar, fiyatlandırma hataları olarak 6 temel kategoride sınıflandırılmıştır. Sağlık alanındaki maliyet konusuna ilave olarak, sağlık alanındaki bir diğer önemli husus ise hastaneler, doktorlar ve diğer sağlık kuruluşları tarafından sağlanan hizmetin kalitesidir. Yapılan çalışmaya [7] göre ABD'de senelik 200-400 bin arasında önlenemez ölüm olduğu ve bu sayının günlük olarak ortalama 1000 kişinin hayatına denk geldiği açıklanmıştır. Sağlık alanındaki hizmet kalitesinin artırılması ve doktorlara yardımcı olabilmesi adına, ESK verilerine yapay zeka yöntemleri uygulanarak, hastalara uygun teşhis/televi önerimi, hastaları önceliklendirme gibi konularda çalışmalar yapılmaya başlanmıştır.

ESK ve diğer sağlık verilerinin artması, dijitalleşmenin yaygınlaşması ve yapay/derin öğrenme yöntemlerinin gelişmesi ile beraber farklı klinik problemleri çözebilmek adına birçok yenilikçi yöntem geliştirilmektedir [8, 9]. Bu araştırmaların genel konuları ise erken hastalık teşhisi ve bu hastalıkları önleme [10], erken mortalite tahmini, hastanede kalma süresi tahmini, hastaneye geri dönme tahmini gibi klinik problemler olup [11], hastalara uygun tedavi önerimi [12], sigorta şirketlerine yönelik çalışmalar [13], ilaç keşfi ve geliştirilmesi [14], epidemiyoloji [15, 16] gibi genel sağlık alanlarında da çeşitli çalışmalar yürütülmektedir. Yapay/derin öğrenme tabanlı çözümler ile sağlık alanındaki problemlere çözümler aranarak hem sağlık alanındaki maliyetleri azaltma hem de daha kaliteli sağlık hizmeti verilmesi hedeflenmektedir. Tez kapsamında, yoğun bakım ünitesinde yatan hastaların ilk 24 içerisinde toplanan klinik verileri farklı çok-kipli (multimodal) derin öğrenme tabanlı yöntemlere girdi olarak verilmiş ve hastanın hastane içerisinde mortalite olma, yoğun bakım ünitesinde mortalite olma, yoğun bakım ünitesinde 3 veya 7 günden fazla kalma ihtimali tahmin edilmeye çalışılmıştır.

1.1 Problemin Tanımı

Derin öğrenme yöntemlerindeki gelişmeler ile birlikte, sayısız olmasına rağmen bazı elektronik sağlık kayıtları setlerinin açık kaynak olarak erişilebilir olması, birçok farklı klinik problemin çözümü için literatürde yeni yöntemler denenmesinin önünü açmıştır. Tez kapsamında yoğun bakımda yatan hastaların ilk 24 saatlik verileri kullanılarak hastaların yoğun bakımda mortalite (in-ICU mortality) ve hastane içerisinde mortalite (in-hospital mortality) olma ihtimalleri ile yoğun bakımda 3 ve 7 günden fazla kalıp kalmayacakları tahmin edilmeye çalışılmıştır. Yoğun bakım ünitesinde hastaların kalma süreleri tahmin edilirken belirlenen 3 ve 7 gün süreleri, üzerinde çalışılan ESK veri seti içerisindeki hastaların yoğun bakımda kalma süre dağılımları ve literatürdeki yoğun bakımda kalma süresini tahmin eden çalışmaların seçtiği gün değerleri dikkate alınarak belirlenmiştir. Yapılan deneylerde hastaya ait farklı özellikler bir arada kullanılarak deneyler çeşitlendirilmiştir. Tez kapsamında çalışılan bu dört klinik problem de ikili sınıflandırma problemi olarak ele alınmıştır. Yapılan deneyler esnasında, hastaya ait yaşamsal gözlem verileri ve laboratuvar sonuçları sabit tutulurken, hastaya ait klinik notların, bu klinik notlar içerisinden çıkartılan medikal terimlerin ve hastaların tedavisi esnasında kullanılan ilaç bilgilerinin çok-kipli derin öğrenme modellerinde kullanılmasının üzerinde çalışılan klinik problemleri tahmin etmedeki etkisi araştırılmıştır.

1.2 Araştırma Motivasyonu ve Önerilen Çözüm Yöntemi

Hastane ve genel sağlık kuruluşlarının temel amacı, sağlık maliyetlerini düşürmek ile beraber sundukları sağlık hizmet kalitesini iyileştirmeye çalışmaktır. Bu sebeple özellikle yoğun bakım ünitelerinde yatan hastaların mortalite oranlarını düşürebilmek hastaneler için önemli bir konudur. Hastanın mortalite olma ihtimalini ilk 24 saatlik yoğun bakım verisinden tahmin ederek, hastaya daha doğru ve planlı tedavi sunabilme imkanı yaratmak insan hayatı için oldukça önemli bir konudur. Diğer yandan bu durum, hastanelerin de ilgili hastalar için daha doğru bir tedavi planı sunmasına olanak sağlamakta ve bu sayede hastane maliyetlerini azaltmasına imkan tanımaktadır. Üzerinde çalışılan diğer klinik problem olan hastaların yoğun bakımda kalma süresinin tahmini de, hem hasta deneyimi hem de yoğun bakım ünitelerindeki yatak ve diğer operasyonların planlanması için kritik bir konudur. Yoğun bakımda yatan hastaların ne kadar süre yoğun bakımda kalacağı kestirimi, yatak planlaması, doktor ve tedavi planlaması gibi birçok konuyu beraberinde etkileyen durumdur. Literatürde mortalite ve yoğun bakımda kalma süresinin tahminine yönelik çalışmalar olmasına rağmen bu çalışmalarda genellikle hastanın sadece zamana bağlı özelliklerinin yani yaşamsal gözlem verilerinin veya laboratuvar sonuçlarının kullanıldığı gözlemlenmiştir. ESK veri seti içerisinde çok farklı veri türleri olduğu göz önüne alındığında bu veri türlerinin

efektif bir şekilde temsil edilmesinin ve bu verilere uygun derin öğrenme mimarileri önerilmesinin üzerinde çalışılan klinik problemlerin tahminine pozitif etki edeceği düşünülmektedir. Bu hedef doğrultusunda, hastanın laboratuvar ve yaşamsal gözlem verilerinin yanısıra, hastaya ait medikal raporlar, doktor ve hemşireler tarafından hasta için yazılan klinik notlar, hastaya verilen ilaçların moleküler bilgileri literatürdeki güncel yöntemler ile temsil edilmiş ve bu temsiller derin öğrenme tabanlı çok-kipli modellere girdi olarak verilerek hastanın mortalite ve yoğun bakımda kalma süreleri tahmin edilmiştir. Sağlık alanı başta olmak üzere diğer birçok alanda oldukça önemli bir konu olan, model tahminlerinin açıklanabilirliği, tez kapsamında ele alınmış bir diğer konudur. Eğitilen modelin mortalite tahminlerini hangi özniteliklere göre yaptığı araştırılmış, ve çıktılar raporlanmıştır. Mortalite tahmini için önemli olan yaşamsal gözlem verileri, laboratuvar deneyleri ve ilaç isimlerinin bulunması, hem derin öğrenme tabanlı modelleri geliştiren kişilerin modellerini test etmeleri için önemli ipucu vermekte, hem de modelin çıktısını kullanacak olan klinik uzmanların model çıktısına daha çok güvenmesini sağlamaktadır.

1.3 Araştırma Zorlukları

Araştırma esnasında karşılaşılan zorlukların önemli bir kısmı sağlık verisinin erişilebilirliği ve heterojen yapısından kaynaklanmıştır. Literatürde açık kaynak elektronik sağlık kayıtları seti sayısı oldukça kısıtlı sayıdadır. Mevcut ESK veri setlerine erişmek için de genellikle çeşitli başvurular ve prosedürlerden geçmek gerekmektedir. Bu durum çalışmaların hızını ve çalışmaya başlama sürecini olumsuz yönde etkileyen önemli bir durumdur. Ayrıca literatürdeki çeşitli çalışmaların özel (açık kaynak olmayan) ESK veri setleri ile çalışması, bu çalışmaların sonuçlarının tekrarlanamamasına ve dolayısıyla bu tez kapsamında önerilen yöntemlerin bu çalışmalar ile karşılaştırılmamasına sebep olmuştur. İkinci olarak, ESK veri setinin içerisindeki verilerin çok çeşitli olması (variety) tasarlanacak olan yapay/derin öğrenme modellerinin yapısını karmaşıktırılmaktadır. Örneğin, hastaya ait zaman serisi tabanlı yaşamsal gözlem verileri ile düzensiz yapıda olan hastaya ait klinik notlar beraber kullanılmak istenildiğinde, farklı türde girdi alan kompleks çok-kipli derin öğrenme tabanlı modeller kurulması gerekmektedir. Araştırma esnasında karşılaşılan bir diğer zorluk ise ESK veri setleri içerisindeki verilerin doğruluğu (veracity) kısmıdır. Hastane içerisinde ölçüm yapan cihazların marka/model farklılıkları, hastanın farklı zamanlarda tekrar tekrar hastaneye gelmesi, farklı zaman aralıklarında gerçekleştirilen laboratuvar deneyleri gibi durumlar, ESK verisi içerisinde birim tutarsızlıkları, eksik veri oluşması gibi problemlere neden olmaktadır. Bu verilerin düzeltilmesi ve standart hale getirilmesi de araştırma alanındaki zorluklardan biridir. Bir diğer husus ise, sağlık verisi ile yapılan çalışmalar için belirli

bir düzeyde alan bilgisi gerekmektedir. Hastalıklara ait kodların yapısı, ilaç ve hastalık isimleri gibi medikal alana özgü kelime ve jargonları bilmek, öznitelik çıkarma ve üzerinde çalışılan problemi formüle etmede önemli bir faktördür. Son olarak ise tez kapsamında dört farklı klinik problem üzerinde model eğitimi gerçekleştirilmiştir. Bu durum, parametre optimizasyonu yaparak optimal model mimarisini bulmayı zorlaştırmış, ve deney sürelerinin zamanını uzatmasına neden olmuştur. Tez kapsamında yapılan çalışmalarda yukarıda sıralanan temel zorluklara çeşitli çözümler bulunmuş, ve ilerleyen bölümlerde bu çözümlerin detaylarından bahsedilmiştir.

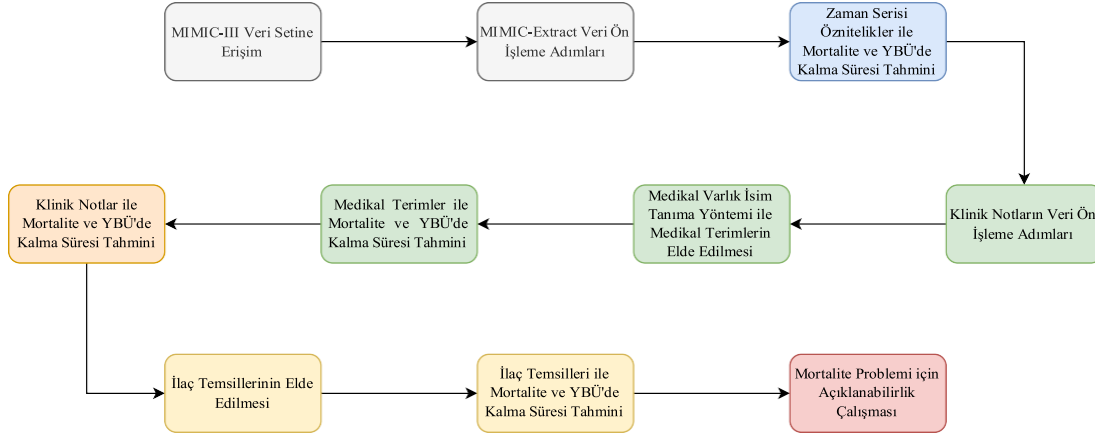
1.4 Tezin Planı

Tez çalışmalarının sırasında izlenen yol bu bölümde okuyucu ile paylaşılmıştır. Hastanın mortalite durumunu ve yoğun bakımda kalma süresini tahmin edebilmek için ilk olarak açık kaynak elektronik sağlık kaydı veri setine erişim sağlanmıştır. İkinci adım olarak bu veri setinin gerekli ön işlemlerden geçirilmesi ve derin öğrenme yöntemlerine uygun girdi formatına dönüştürülmesi süreci gerçekleştirilmiştir. Ardından, elde edilen zaman serisi öznitelikleri ile mortalite ve yoğun bakımda kalma sürelerinin tahminine yönelik modeller eğitilmiştir. Daha sonrasında, modellerin tahmin başarımlarını arttırabilmek için hastalara ait klinik notların modellere girdi olarak verilebilmesi üzerine çalışmalar gerçekleştirilmiştir. Bu kapsamda, klinik notlar veri ön işlemlerinden geçirilmiş, varlık isim tanıma yöntemi ile medikal terimler elde edilmiş, ve zaman serisi öznitelikleri ile medikal terimler beraber kullanılarak modeller tekrar eğitilmiştir. Yapılan diğer çalışmada ise klinik notları içerisinden medikal terimleri çıkartarak kullanmak yerine, doğrudan klinik notları kullanılmış ve model tahminlerine etkisi araştırılmıştır. Yapılan son çalışmada ise, hastaya ait ilaç bilgilerinin moleküler temsillerini zaman serisi öznitelikleri ile beraber kullanmanın avantajları araştırılmıştır. Bunun için ilk olarak ilaçların temsilleri öğrenilmiş ve ardından ilgili modeller eğitilmiştir. Son aşamada ise, mortalite problemi için açıklanabilirlik konusu çalışılmış, ve mortalite tahmini sırasında hastaya ait hangi zaman serisi özniteliklerin ve ilaçların önemli olduğu tespit edilmeye çalışılmıştır. Anlatılan bu adımlar Şekil 1.1'de okuyucu ile paylaşılmıştır.

1.5 Tezin Katkıları

Tez içerisinde mortalite ve yoğun bakımda kalma süresini tahmin etmek için geliştirilmiş üç ayrı yöntem bulunmaktadır. Bu yöntemlerin literatüre katkılarına aşağıda listelenmiştir:

1. Hastaya ait zamana bağlı yaşamsal gözlem verileri ve laboratuvar sonuçlarının yanısıra hastaya ait klinik notlar, dönüştürücü (transformer) tabanlı bir yöntem olan BERT (Bidirectional Encoder Representations from Transformers) modeli



Şekil 1.1: Tezin temel adımları.

ile temsil edilerek, klinik notların mortalite ve yoğun bakımda kalma süresinin tahmini problemlerine etkisi araştırılmıştır. Ayrıca bu temsiller ile zamana bağlı öznitelikleri efektif bir şekilde kullanabilmek için derin öğrenme tabanlı çok-kipli model önerilmiştir.

2. Klinik notların çok uzun ve çeşitli tıbbi jargonlar içermesi sebebiyle içerisindeki bu kirli ve gürültülü veriyi temizlemek için medikal terimlerin çıkartılması hedeflenmiştir. Klinik notlar içerisinde medikal terimlerin çıkartılabilmesi için klinik alana özel eğitilmiş varlık isim tanıma yöntemi kullanılmıştır. Çıkarılan medikal terimler zaman-serisi öznitelikleri ile beraber kullanılarak mortalite ve yoğun bakımda kalma süresi tahmin edilmiştir. Bu tahminler, derin öğrenme tabanlı çok-kipli model ile gerçekleştirilmiştir. Model başarımlarının iyileştirilmeye çalışılmasına ek olarak, deneyler kapsamında klinik notlar içerisinde çıkarılan medikal terimlerin vektörel temsillerinin başarılı bir şekilde yapılabilmesi için farklı kelime gömme (word embedding) yöntemleri uygulanmıştır.
3. Yoğun bakım ünitesinde yatan hastalara uygulanan tedavi kapsamında verilen ilaçların bilgisinin, mortalite ve yoğun bakımda kalma sürelerini tahmin etmede, zaman serisi öznitelikleriyle beraber kullanılması önerilmiştir. İlaçların moleküler bilgilerini vektörel hale dönüştürebilmek için hem geleneksel moleküler parmak izi çıkartma yöntemleri hem de yapay öğrenme tabanlı moleküler temsil öğrenme çalışmalarından yararlanılmıştır. Çok-kipli derin öğrenme tabanlı modellerde, ilaçların moleküler temsillerinin kullanılmasının, mortalite ve yoğun bakımda kalma süresini tahmin etme problemlerinde kullanılması, bildiğimiz kadarıyla literatürde ilk kez denenmiştir.
4. Mortalite tahmini yapmak için, hastaya ait yaşamsal gözlem verilerinin, laboratu-

var sonuçlarının ve ilaçların bilgilerini kullanarak eğitilen derin öğrenme tabanlı modelin açıklanabilirliğini arttırabilmek adına yöntem önerilmiştir. Bu sayede mortalite tahmini için önemli olabilecek yaş, amsal gözlem verileri, laboratuvar deneyleri ve önemli klinik ilaçlar tespit edilmeye çalışılmış, ve okuyucular ile paylaşılmıştır.

1.6 Tezin Düzeni

Bölüm 2’de ESK verileri başta olmak üzere, birçok sağlık verisi üzerine gerçekleştirilen yapay öğrenme yöntemlerinden ve klinik problemleri çözmeye çalışan literatürdeki diğer güncel çalışmalardan bahsedilmiştir. Tez kapsamında gerçekleştirilen deneyler esnasında kullanılan veri seti ve bu veri seti üzerinde uygulanan veri ön işleme adımları ise Bölüm 3’te anlatılmıştır. Bölüm 4’de, yapılan deneyler esnasında kullanılan yapay öğrenme ve derin öğrenme yöntemleri, kelime temsil yöntemleri, açıklanabilirlik yöntemi, ilaç temsilleri gibi algoritmaların detaylarından bahsedilmiştir. Bölüm 5’de, klinik notları doğrudan kullanmak yerine, içerisinden medikal terimlerin çıkartılarak zaman serisi verileri ile beraber kullanılmasının detayları ve sonuçları anlatılmıştır. Bölüm 6’te ise üzerinde çalışılan klinik problemlerin çözümü için klinik notların nasıl temsil edildiği ve zaman serisi verileri ile beraber nasıl kullanıldığının anlatılan yöntem tartışılmıştır. Bölüm 7’de hastaya ait ilaç bilgilerini model içerisinde kullanmanın, klinik problemlere etkisi araştırılmış, ve sonuçlar okuyucu ile paylaşılmıştır. Son olarak ise, Bölüm 8’de konuyla ilgili toparlayıcı sonuç bilgiler ve yapılan araştırmanın ileride nasıl devam edebileceği ile ilgili öneriler verilmiştir.



2. İLGİLİ ÇALIŞMALAR

Derin öğrenme yöntemleri, son yıllarda araştırmacıların görüntü işleme, ses işleme, doğal dil işleme gibi farklılandaki birçok zorlu problem üzerinde en iyi sonuçları almasına olanak sağlamıştır [17]. Farklı problemler üzerinde derin öğrenme yöntemleri sayesinde alınan başarılı sonuçlar ile sağlık alanında da bu yöntemler uygulanmaya başlanmıştır. Klinik veriler, genellikle zaman damgasıyla kaydedilmekte ve bir hastaya ait veriler zaman serisi verisi olarak ele alınabilmektedir. Zaman serisi tabanlı bu verileri işleyebilmek ve çeşitli klinik konseptlerin temsillerini öğrenebilmek adına literatürde çalışmalar gerçekleştirilmiştir. Bu alandaki popüler ve ilk çalışmalardan sayılan Med2Vec [18] yöntemi 2016 yılında Choi vd. tarafından önerilmiştir. Çalışmanın amacı hastaya ait hastane ziyaretlerinin ve medikal kodların temsillerini öğrenmeye yönelik yapay öğrenme tabanlı bir mimari geliştirmektir. Çalışmada, Children's Healthcare of Atlanta (CHOA) hastanesi verisi kullanılmıştır. CHOA veri seti, içerisinde 550,339 tekil hasta, 3,359,240 tekil hastane ziyareti, hasta başına ortalama 6.1 hastane ziyareti ve her ziyarette ortalama 7.88 medikal kod bulunduran büyük bir ESK veri setidir. Girdi olarak kullanılan medikal kodlar kendi içerisinde üçe ayrılmaktadırlar: teşhis, ilaç (tedavi), ve prosedür kodları. Önerilen yöntemde, 2 katmanlı (medikal kodlar ve hastane ziyaretlerini öğrenebilmek için ayrı ayrı) bir yapay sinir ağı mimarisi önerilmiştir. Önerilen mimariye girdi olarak her bir hastane ziyareti içerisinde yer alan medikal kodların vektörel temsilleri verilmiştir. Ardından, girdi iki lineer katmanlı yapay sinir ağından geçirilerek önceki ve sonraki hasta ziyaret vektörlerini tahmin etmeye çalışan bir mimari önerilmiştir. Öğrenilen medikal kod temsilleri, hastanın gelecek ziyaretlerindeki medikal kodlarını ve hasta risk gruplarını tahmin etmek için kullanılmıştır. Elde edilen sonuçlar, literatürdeki Word2Vec [19] (skip-gram), Global Vectors for Word Representations (Glove) [20], yığın otomatik kodlayıcı (stacked autoencoder) gibi yöntemlerle karşılaştırılmış ve okuyucu ile bu yöntemlere göre

daha başarılı sonuçlar alındığı raporlanmıştır. Bu alandaki bir diğer çalışma olan MIME [21], 2018 yılında Choi vd. tarafından yayınlanmıştır. Med2Vec çalışması ile benzer motivasyonu olan MIME çalışmasında, temsil öğrenme mimarisi geliştirilerek, tedavi, teşhis, ziyaret ve hasta seviyesindeki kodları ayrı ayrı öğrenmeye olanak sağlayan bir mimari önerilmiştir. Sutter Health kurumunun verisi kullanılarak eğitilen bu model,

30,764 hasta ve 616,073 hastane ziyaret verisi ile eğitilmiştir. Öğrenilen temsillerin başarımının gösterilebilmesi adına, yine aynı veri seti kullanılarak hastaların kalp krizi riski tahmini yapılmıştır. Deney sonuçları incelendiğinde, önerilen modelin Med2Vec, GRAM ve diğer yöntemlerden daha başarılı sonuçlar verdiği gösterilmiştir. GRAM [22], medikal kodların temsilini dikkat ağırları (attention) üzerinden çizge tabanlı bir yöntem ile öğrenmeyi önermektedir. Özellikle az miktarda veri ile yapılan deneylerde, karşılaştırıldığı diğer çalışmalara göre daha yüksek oranda Area Under the Receiver Operating Characteristic Curve (AUROC) skoru verdiği gösterilmiştir. Bir diğer çalışma olan Deep Patient [23] isimli çalışmada ise Mount Sinai Hastanesindeki 700 bin hasta verisi kullanılarak deneyler gerçekleştirilmiştir. Hastalara ait birçok veri türü (hastalık, tedavi, prosedür kodları, lab testleri, klinik notlar, demografik bilgileri) çok katmanlı gürültü arındırıcı otokodlayıcılara (multiple denoising autoencoders) girdi olarak verilerek hastanın vektörel temsili öğrenilmeye çalışılmıştır. Öğrenilen hasta temsilleri, rastgele orman (random forest) algoritmasına girdi olarak verilerek 78 farklı hastalık kodu tahmin edilmiştir.

Literatürdeki yöntemlerin büyük çoğunluğunun nihai amacı hastaya ait çeşitli çıktıları tahmin etmektir. Choi vd. [10] yaptığı çalışmada Tekrarlayan Sinir Ağları (Recurrent Neural Network, RNN) [24] modelini kullanılarak hastanın kalp krizi geçirme riskini tahmin etmiştir. Tahmin esnasında hastanın 12-18 aylık geçmiş, teşhis, tedavi ve prosedür bilgileri kullanılmıştır. Sutter Health kurum verisi içerisinde yer alan ve 4 bini vaka olmak üzere 34 bin hasta üzerinde gerçekleştirilen çalışmada, 34 bin hastanın seçiminde hastaların birden çok kez hastaneye gelmiş olmalarına dikkat edilmiş, ayrıca hastalar 40-85 yaş arasından seçilmiştir. Sonuçlar incelendiğinde 0.83'lük bir AUROC skoru elde edildiği raporlanmıştır. DoctorAI [25] ismi verilen derin öğrenme tabanlı model ise hastalardaki hastalığın ilerleyişini modellemek için önerilmiştir. Hastanın geçmiş teşhis bilgileri kullanılarak bir sonraki hastane ziyaretinde hastaya hangi teşhisin konulabileceği tahmin edilmiştir. Deneyler, Sutter Health kurum verisinden 260 bin hastanın 10 yıllık verisi üzerinde gerçekleştirilmiştir. Deneyler sonucunda önerilen derin öğrenme tabanlı RNN modelinin diğer temel yöntemleri geçtiği bilgisi okuyucular ile paylaşılmıştır. RETAIN [26] çalışması, 263 bin hastanın yaklaşık 14 milyon hastane ziyaret verisini kullanarak hastanın kalp rahatsızlığı geçirip geçirmeyeceğini tahminleyen bir model önermektedir. Önerilen bu model RNN tabanlı olmasının yanı sıra aynı zamanda dikkat mekanizması (attention mechanism) içermekte olup, bu sayede açıklanabilir tahminler üretebilmektedir. Çalışma kapsamında önerilen model 0.87 AUROC skoru ile tahmin edilirken, yapmış olduğu bu tahminleri hastanın hangi hastane ziyaretine ve bu ziyaretlerdeki hangi hastalık teşhisine ne kadar önem verdiği bilgisini de açıklayabilmektedir. Tahminlerin açıklanabilir olması, modeli kullanacak doktor ve

diğer klinik çalıřanlar için oldukça önemli bir özelliktir.

Literatürde, hastalıkların gelecekteki durumunu tahmin etmeye yönelik birçok sayıda çalıřma mevcut olmasıyla beraber [27, 28, 29], farklı klinik problemlerde çözülmeye çalışılmaktadır. Örneğin, hastanın yaptığı ilk hastane ziyaretinden sonra 30 gün içerisinde hastaneye tekrar gelmesi (readmission) önemli bir konudur. Fatuma vd. [30] çalışmasında da bahsedildiği gibi hastane başvurularının %17'si hastanın 30 gün içerisinde hastaneye yaptığı tekrar ziyaretlerini kapsamaktadır. Bu %17'lik kısmın ise %75'inin aslında önlenabilir olduğu ve bunun çok büyük maliyetlere sebep olduğu ortaya konmuştur. Bu sebeple New Zeland National Minimum Veri seti ile yapılan çalışmada hastaların hastaneyi tekrar ziyaret edip etmeyeceği hastanın geçmiş verileri kullanılarak tahmin edilmeye çalışılmış ve önerilen derin yapay sinir ağları diğer yöntemlerden daha başarılı sonuç göstererek %73.4 AUROC skorunu vermiştir. Bu çalışmalar haricinde literatürde oldukça farklı klinik problemler üzerinde de çalışmalar oldukça yoğun bir şekilde devam etmektedir. Klinik notlar içerisinde International Classification of Disease (ICD) kodlarının çıkartılması [31], sentetik ESK verisi üretme [32], sentetik hasta verisi üretme [33], hastaya ait görüntülerden deri kanseri [34] tahmini, diyabetik retinopati tespiti [35] gibi birçok farklı alanda çalışma gerçekleştirilmektedir. Hastaya ait geçmiş, hastalık bilgileri, yaşamsal gözlem verileri gibi veri tipleri haricinde, hastaya ait klinik notların doğal dil işleme yöntemleri ile işlenerek üzerinde çalışılan klinik problem için girdi olarak sağlanmasında literatürde çalışılmaya başlanan bir araştırma konusudur. Doğal dil işleme alanında 2013 yılından itibaren önemli yöntemler yayınlanmıştır. Word2vec [19], Glove [20], FastText [36] gibi temsil öğrenme yöntemleri önerilmiş ve kelimelerin uzayda daha iyi bir şekilde temsil edilmesi sağlanmaya çalışılmıştır. Ayrıca kelime yerine doğrudan bir dokümanın temsilini öğrenmeyi öneren Doc2Vec [37] çalışması yayınlanmıştır. Bu gelişmelerle beraber normalde görüntü üzerinde iyi çalıştığı bilinen Evrişimsel Sinir Ağları (CNN) metin üzerinde uygulayan, bu sayede hızlı ve başarılı sonuçlar elde eden 1D Evrişimsel Sinir Ağları [38] da metinler üzerinden öznetelik çıkarımı için kullanılmaya başlanmıştır. Bu ve benzeri gelişmelerin ortaya çıkması ile elektronik sağlık kayıtları içerisindeki klinik notların problem çözümü için önemi artmış olup çeşitli çalışmalarda kullanılmaya başlanmıştır. Liu vd. [39], zaman serisi özneteliklerinin yanısıra yapısal olmayan klinik notları da kullanılarak kronik rahatsızlıkların tahminini gerçekleştirmiş ve klinik notları kullanmanın model performansına olumlu etkisini göstermiştir. Si ve Roberts [40] ise klinik notları kullanarak hastanın hastanede ve hastaneden çıktıktan 30 gün veya 1 sene içerisinde mortalite olma ihtimalini tahmin etmiştir. Boag vd. [41], Bag of Words (BoW), Word2Vec gibi kelime temsil yöntemlerini karşılaştırarak klinik notların daha iyi nasıl temsil edileceğini araştırmış ve modellerin başarımını hastalık ve mortalite tahmini

problemleri üzerinde test etmiş, tir. Bir başka çalışmada ise klinik notlar içerisinde medikal kodların ortaya çıkartılması üzerine çalışma gerçekleştirilmiş, tir [42]. Çok sınıflı bir sınıflandırma problemi olarak tasarlanan bu problemde bu tezde de kullanılan Medical Information Mart for Intensive Care (MIMIC-III) [43] veri seti kullanılmış, ve evrimsel sinir ağı ile birlikte dikkat mekanizması kullanılarak derin öğrenme tabanlı yöntem önerilmiş, tir.

Bu gelişmelere ek olarak doğal dil işleme alanında yakın tarihte çok önemli bir gelişme daha olmuştur. 2018 yılında Google'ın yayınlamış olduğu Bidirectional Encoder Representations from Transformers (BERT) [44] modeli ile doğal dil işleme alanındaki birçok problem üzerinde en iyi sonuçlar alınmaya başlanmıştır. Bu gelişmelerle birlikte elektronik sağlık kayıtları içerisindeki klinik notlar üzerinde BERT ve benzeri dönüştürücü [45] (transformer) tabanlı yöntemlerle klinik alanda yeni temsiller öğrenilmiştir. Alsentzer vd. [46] MIMIC-III veri seti içerisindeki klinik notları kullanarak BERT ve BioBERT [47] modelleri üzerinde ince ayar yaparak klinik temsilleri öğrenmiş ve bu öğrenilen klinik temsilleri araştırmacılar ile paylaşmışlardır. Huang vd. [48] ise çalışmalarında bir önce açıklanan çalışmayla benzer bir şekilde MIMIC-III veri seti içerisindeki klinik notları kullanarak BERT modelini ince ayar ile klinik alana özgü olarak eğitmiştir. [46] çalışmasından farklı olarak ise ince ayar esnasında eğitimi yapılan klinik problem hastaneye tekrar başvuru problemi olarak seçilmiş ve eğitilen klinik alanına özgü BERT modeli herkese açık hale getirilerek paylaşılmıştır. Literatürde klinik notlar ile yapılan çalışmalar olmasına rağmen, klinik notların uzunluğu, doktorların klinik notları yazarken medikal jargon kullanmaları, kısaltma ifadelerine yer vermeleri gibi sebeplerden dolayı klinik notları ön işleme tabi tutmak ve modele girdi olarak vermek zorlu bir süreçtir. Bu sebeple çeşitli çalışmalarda da klinik notları doğrudan kullanmak yerine içerisindeki medikal terimler çıkartılarak bu terimlerin kullanılması önerilmiştir [49, 50, 51]. Scispacy [52] kütüphanesi klinik/biyomedikal alandaki metinler üzerinde kelime türü etiketleme (pos tagging), bağımlılık analizi (dependency parsing), varlık isim tanıma (name entity recognition) gibi yöntemleri kullanmayı sağlamaktadır. Bu tez kapsamında ise klinik notlar içerisindeki medikal terimlerin çıkartılabilmesi için med7 [53] çalışması kullanılmıştır. Med7 çalışmasında, model öncelikli olarak bir sonraki kelimeyi tahmin etme problemi üzerinde öz-denetimli öğrenme (self-supervised learning) ile eğitilmiş daha sonra ise MIMIC-III içerisindeki az sayıdaki etiketli klinik notlar ile varlık isim tanıma problemi ile ince ayarı yapılmıştır. Eğitilen bu klinik alandaki varlık isim tanıma modeli ilaçların ismi, dozu, gücü, formu, sıklığı, süresi, verilme yöntemi olmak üzere 7 farklı varlık ismi çıkartılabilmektedir. Klinik alandaki derin öğrenme yöntemleri ile doğal dil işleme çalışmaları hakkında daha detaylı bilgi için Wu vd. [54] tarafından yapılmış olan literatür taraması okuyucular tarafından incelenebilir.

ESK içerisindeki bir başka veri türü ise hastalara verilen ilaç bilgileridir. İlaçların farklı temsilleri, keminformatik (cheminformatic) uygulamaları içerisinde yer alan ilaç keşfi [55], ilaç-ilaç etkileşim tahminleri [56], bileşik-protein yakınlığı tahmin etme [57] gibi çalışmalarında yoğun olarak kullanılmaktadır. İlaç bilgilerini yapay öğrenme ve derin öğrenme problemlerine girdi olarak verebilmek için ilaçların temsillerini öğrenmeyi öneren farklı yöntemler bulunmaktadır. Bu tez kapsamında Extended-Connectivity Fingerprints (ECFP) [58], Molecular ACCess System (MACCS) [59], Mol2Vec [60], Smiles-Transformer [61] yöntemleri ile deneyler gerçekleştirilmiştir. Bu temsil yöntemleri ile ilgili detaylar Bölüm 4.4'de okuyucu ile detaylı olarak paylaşılmıştır. Derin öğrenme alanında farklı veri türleri beraber kullanılarak çok-kipli (multimodal) yöntemler ile görüntülerden açıklama oluşturma [62] (image captioning), görüntüler üzerinden soru cevaplama [63] (visual question answering), ses tanıma [64] gibi farklı problem türlerinde başarılı sonuçlar alınmaktadır. Sağlık alanında da hastaya ait zamana-bağlı yaşamsal gözlem verileri, klinik notlar, ilaç temsilleri gibi farklı veri türleri beraber kullanılarak üzerinde çalışılan klinik problemler üzerinde daha başarılı sonuçlar alınması hedeflenmiştir. Khadanga vd. [65] tarafından yapılan çalışmada yapısal olmayan klinik notlar ile yapısal zaman-serisi öznitelikleri birlikte kullanılarak hastane içerisindeki mortalite, organ yetmezliği ve hastanede kalma süresi tahmin edilmeye çalışılmıştır. Benzer bir şekilde, [66] çalışmasında ise mortalite tahmininde fizyolojik zaman serisi verileri ve klinik notların nasıl entegre edilebileceği keşfedilmeye çalışılmıştır. Jin vd. [67] yaptığı çalışmada klinik notları doğrudan kullanmak yerine, notlar içerisindeki medikal terimlerin çıkartılmasını önermiştir. Çıkartılan medikal terimler Doc2Vec [68] yöntemi ile temsil edilirken, hastaya ait zaman serisi verileri Uzun Kısa Vadeli Hafıza Ağları (Long Short-Term Memory, LSTM) [69] ile temsil edilmiş ve bu temsiller birleştirilerek hastanın mortalite olma ihtimali tahmin edilmiştir. Bu tez kapsamında da geliştirilen çalışmalar zaman-serisi öznitelikleri ile beraber klinik not [70], medikal terim [71], ilaç bilgisi [72] gibi farklı veri türlerini bir arada kullanarak çok-kipli derin öğrenme tabanlı yöntemler önerilmiştir. Çok-kipli yöntemler ile yapılan çalışmaların geniş bir literatür taraması için Baltrusaitis vd. [73] yapmış olduğu çalışma incelenebilir.

Son zamanlarda yapay ve derin öğrenme algoritmaları birçok alanda etkileyici sonuçlar göstermesine ek olarak, modellerin yaptıkları tahminleri neden yaptığina dair anlamlandırma işlemi de oldukça önemli bir problemdir. Modellerin açıklanabilirliği, bu yöntemleri kullanan kişilerin modellere olan güvenini artırırken aynı zamanda modelin güvenen kişiler için de modelin içerisindeki olası problemleri keşfetmesini sağlamaktadır [74]. Lineer regresyon, lojistik regresyon, karar ağaçları gibi kendiliğinden yorumlanabilen algoritmalar olduğu gibi, çoğu yöntemi açıklanabilir hale getirmek veya tahminlerini

anlamlandırabilmek için ek yöntemlere ihtiyaç duyulmaktadır. Bu tez kapsamında son zamanlarda oldukça popüler olan SHapley Additive exPlanations (SHAP) [75] yöntemi kullanılmaktadır. Bu yöntem ile siber güvenlik [76], trafik kazası tahmini [77], sağlık uygulamaları [78] gibi farklı alanlarda çalışmalar gerçekleştirilmiştir. Yakın zamanda ise, Rodriguez-Perez ve Bajorath tarafından önerilen çalışmalarda [79, 80] bileşiklerin etki gücü ve çok hedefli aktivite tahmini için önemli olan moleküler alt yapılarını açıklayabilmek için SHAP yönteminden faydalanılmıştır. Sağlık alanında yapay zeka tabanlı yöntemler kullanılarak gerçekleştirilen çalışmalar, bu çalışmalar yapılırken karşılaşılan zorluklar ve gelecekte bu alanda yapılabilecek fırsatlar ile ilgili daha detaylı ve sistematik bilgi için literatürdeki [81, 82, 4, 83] çalışmalar incelenebilir.



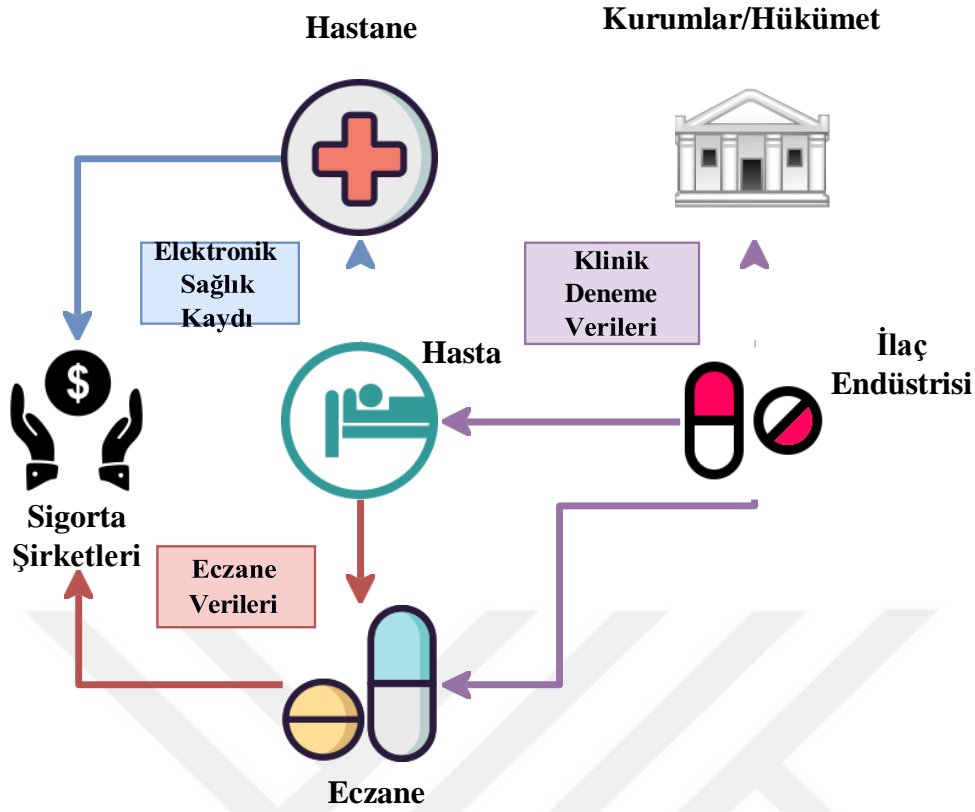
3. VERİ KÜMESİ

Sağlık verisi, birçok kurum ve kuruluşun beraber ürettiği bir veri türüdür. Merkezinde hastanın geçmişine ait verilerin olduğu sağlık verisinde, hastaların hastaneye gitmesi, hastanelerle sigorta şirketlerinin bilgi alışverişi ile beraber sağlık verisinin bir alt kümesi olan ESK veri setleri oluşmaktadır. Buna ek olarak, hastaların hastane ziyaretinden sonraki eczaneye ziyaretleri, eczanelerin sigorta şirketleri ile etkileşim kurmaları ve ilaç alım/satım işlemi ile birlikte eczane talep verileri oluşmaktadır. Bu veride hastalık-teşhis-ilaç arasında bir veri seti oluşturularak farklı araştırma alanları yaratılmaktadır. Sağlık verisi altındaki bir diğer veri türü ise klinik denemelerin sonucunda oluşmaktadır. İlaç firmalarının, hastalar, çeşitli kurumlar ve hastaneler ile iş birliği yaparak gerçekleştirmiş oldukları klinik denemelerin verisi de sağlık verisi altında sınıflandırılmaktadır. Bahsi geçen tüm bu verileri kullanarak literatürde araştırmalar ve yayınlar yapılmaktadır. Bu anlatılan sağlık verisinin döngüsü Şekil 3.1 paylaşılmıştır.

Tez kapsamında gerçekleştirilen deneylerde kullanılmak için ESK veri kümesine ihtiyaç duyulmuştur. Açık kaynaklarda bulunan ESK veri setlerinin sayısı, sağlık verisinin mahremiyeti sebebiyle oldukça az sayıdadır. Bu tezde, literatürde en çok kullanılan açık kaynak ESK veri seti olan Medical Information Mart for Intensive Care (MIMIC-III) [43] kullanılmıştır. Bu bölümde ise, kullanılan veri setinin detayları Bölüm 3.1 ve veri setinin yapılan farklı deneyler için nasıl ve hangi işlemlerden geçirilerek özniteliklere dönüştürüldüğü Bölüm 3.2 açıklanmıştır.

3.1 MIMIC-III Veri Kümesi

Gerçekleştirilen tüm deneylerde, ESK veri seti olarak Medical Information Mart for Intensive Care (MIMIC-III) [43] kullanılmıştır. Bu veri seti, içerisinde 2001-2012 yılları arasında Beth Israel Deaconess Medical Center (Boston, Massachusetts) hastanesinin çeşitli yoğun bakım ünitelerinde kalmış olan 46,520 hasta, tekil 58,976 hastane ziyareti ve 61,532 yoğun bakım ziyareti verisini içermektedir. [43] çalışmasından alınılanarak çizilen Şekil 3.2’de görüleceği üzere, hastanenin farklı birimlerinden toplanan veriler önce birleştirilmiş, ardından kimliği gizleme, tarihleri kaydırma ve format düzenlemesi yapılarak hem veri anonimleştirilmiş, hem de araştırmacılar için kullanıma hazır hale getirilmiştir. Veri seti, hastanın demografik bilgileri, hastaneye giriş/çıkış bilgileri,



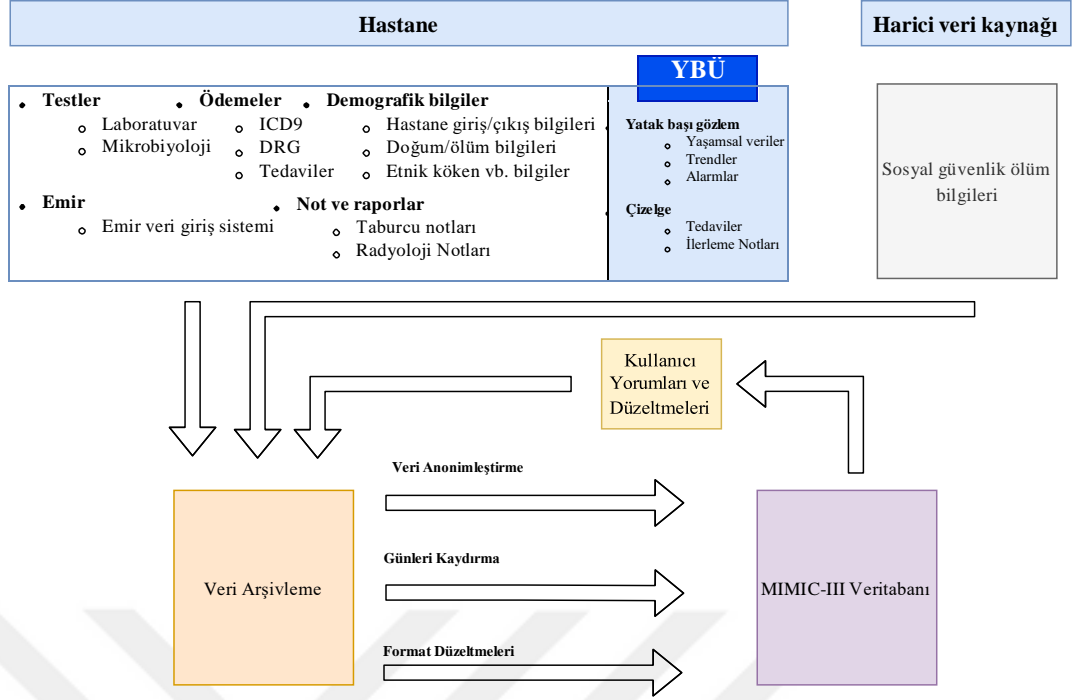
Şekil 3.1: Çeşitli sağlık verilerinin oluşum şekilleri.

hastaya uygulanan tedaviler, klinik notlar, laboratuvar sonuçları, yatak başı gözlem verileri gibi birçok farklı veri türünü içermektedir. Veriler ve verilerin tutulduğu tablolar ile ilgili detay bilgiler ise Bölüm 3.1.2’de okuyucu ile paylaşılmıştır.

3.1.1 MIMIC-III Veri setine erişim

MIMIC-III veri setine erişmek için tamamlanması gereken çeşitli kurslar ve araştırmacı tarafından imzalanması gereken veri kullanım şartı bulunmaktadır. Bunun için öncelikli olarak Collaborative Institutional Training Initiative (CITI) kurumu tarafından hazırlanan "Data or Specimens Only Research" kursunun tamamlanması gerekmektedir. Bu kurs 9 ayrı modülden oluşmakla beraber, MIMIC-III verisini kullanacak olan araştırmacıya, sağlık verilerinin hassasiyeti, etik kurallar, Health Insurance Portability and Accountability Act (HIPAA) gereksinimleri ve benzeri konular ile ilgili bilgi vermeyi amaçlamaktadır. Alınan bu kurslardan sonra yapılan sınavı başarıyla geçen araştırmacılar, bu kursu tamamlamış sayılmaktadırlar. Kursu başarıyla tamamlayan araştırmacıların PhysioNet¹ üzerinden hesap ve MIMIC-III veri seti için talep açmaları gerekmektedir. Erişim onayılan ve sağlık verisini kullanma talimatlarını içeren belgeyi

¹<https://physionet.org/pnw/login>



Şekil 3.2: MIMIC-III genel mimarisi.

imzalayan araştırmacılar, MIMIC-III veri setine erişim sağlayabileceklerdir. Yaklaşık 40GB büyüklüğünde, 26 farklı "CSV" dokümanından oluşan MIMIC-III veri seti, lokal ortamda doğrudan CSV doküman üzerinde veya ilişkisel veri tabanına aktarılarak kullanılabilir. Aynı zamanda bulut ortamda çalışmaya yapmak isteyen araştırmacılar için Google BigQuery, Amazon Web Service (AWS), Google Cloud Storage (GCS) seçenekleri mevcuttur.

3.1.2 MIMIC-III Tablo detayları

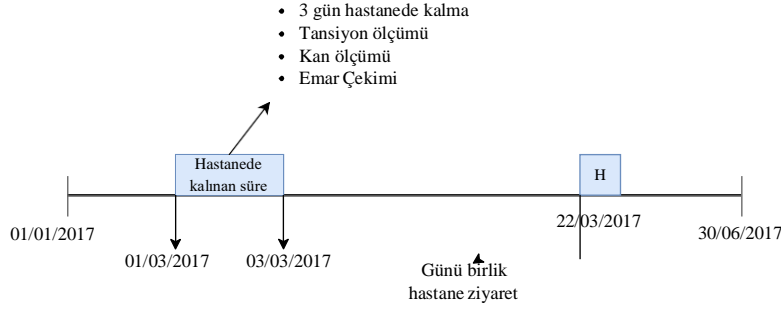
MIMIC-III veri seti içerisinde hastaya ait demografik bilgiler, yaşamsal gözlem verileri, klinik ölçümler, ödeme bilgileri, hastanın medikal geçmişi, uygulanan tedaviler vb. diğer tüm veriler 26 farklı tablo içerisinde saklanmaktadır. Farklı tablolar içerisindeki aynı hastaya, hasta ziyaretine veya yapılan bir deneye ait sonuçları elde edebilmek için ise tablolar içerisinde birincil anahtarlar mevcuttur. Bu eşleştirme işlemi tablolar içerisindeki son eki 'ID' olan sütun isimleri ile gerçekleştirilebilmektedir. Örneğin, SUBJECT_ID her bir tekil hastayı, HADM_ID her bir tekil hastane ziyaretini ve ICU_ID ise her bir tekil yoğun bakım ziyaretini ifade etmektedir. Yoğun bakımda kalan hastalardan toplanan veriler ise sonu 'events' ifadesi ile biten tablolarda tutulmaktadır (örnek: CHARTEVENTS, NOTEEVENTS, LABEVENTS vb.). Ön ek olarak 'D_' ifadesi ile başlayan tablolar ise sözlük görevi görmekte olup, diğer tablolarda yer alan kodların açıklamalarını içermektedir (örnek: D_ITEMS).

MIMIC-III içerisindeki tablo yapısı özetle şu şekilde açıklanabilir. Hastaların tanımlamak ve hastane içerisindeki durumunu takip etmek için 5 tablo bulunmaktadır: ADMISSIONS, PATIENTS, ICUSTAYS, SERVICES ve TRANSFERS. Tablolar arasındaki bağlantıları sağlayan ve çeşitli terimlerin açıklamalarını içeren de 5 ayrı tablo bulunmaktadır: D_ICD_DIAGNOSIS, D_ICD_PROCEDURES, D_ITEMS, D_LABITEMS, D_CPT. Bunların haricinde kalan diğer 16 tablo ise hastanın durumu, fiziksel ölçümleri, yaşamsal belirtileri, klinik notları, tıbbi görevli bilgileri gibi diğer bilgileri içermektedir (hastane ve yoğun bakım ünitesi verileri).

Yapılan deneylerde doğrudan kullanılan tabloların içerikleri daha detaylı olarak tartışılmıştır. MIMIC-III içerisindeki en temel tablolardan biri olan ADMISSIONS tablosu, hastaneye yapılan her bir giriş, i (bas, vuru) tekil olarak saklanmaktadır. Tablo içerisinde, hastanın hastaneye giriş, /çıkış, tarihleri, bas, vuru tipi, bas, vuru yeri, demografik bilgileri (din, medeni durum, etnisite vb.), bas, vuru ş, ikayeti, sigorta bilgisi, mortalite durumu gibi bilgiler bulunmaktadır. Her bir hastane başvurusu, tekil olarak HADM_ID değişkeni altında saklanmakta ve bir hastanın tekil hastane ziyareti diğer tablolar ile HADM_ID değişkeni üzerinden kurulmaktadır. ADMISSIONS tablosu ile bağlantılı bir diğer temel tablo ise PATIENTS tablosudur. PATIENTS tablosu, her bir hastanın cinsiyet, doğum tarihi, eğer mortalite olduysa mortalite tarihini tutmaktadır. Her bir hasta ise tekil olarak SUBJECT_ID değişkeni ile tanımlanmaktadır. ICU_STAYS tablosu, yoğun bakım ünitesindeki tekil ziyaret bilgilerini içermektedir. Bu tablo içerisinde, hastanın yatmış, olduğu ilk ve son yoğun bakım ünite tipleri, yoğun bakıma giriş ve çıkış,

zamanları, yoğun bakım ünitesinde kaldığı toplam süre gibi bilgiler bulunmaktadır. Her bir yoğun bakımda kalma, ICUSTAY_ID değişkeni ile tekil olarak temsil edilmekte ve her bir ICUSTAY_ID değişkeninin bağlı olduğu bir SUBJECT_ID ve HADM_ID olmak zorundadır. Örnek bir ESK toplanma akışı Şekil 3.3'de gösterilmiştir.

Hastanın YBÜ'de kaldığı süre boyunca toplanan verileri CHARTEVENTS tablosunda saklanmaktadır. Bu veriler, baş, ta hastanın rutin yaş, amsal belirtileri ve laboratuvar sonuçları olmak üzere elektronik ekranlarda gösterilen diğer tüm değerleri taşımaktadır. Tablo içerisinde bulunan ITEMID alanı, tekil olarak yapılan her bir ölçümün tanımlayıcısıdır. Bunun haricinde bu ölçümün ne zaman yapıldığına zaman saklandığı, kim tarafından yapıldığı ve ölçüm değerleri gibi bilgiler yer almaktadır. Tez kapsamında gerçekleştirilen deneylerde kullanılan hastaya ait klinik notlar ise NOTEEVENTS tablosunda bulunmaktadır. Klinik notlar, doktor ve hemş, irelerin hastalar için yazmış, oldukları rutin gözlemler ile ilgili olabilece ği gibi, radyoloji, elektrokardiyogram gibi farklı konular ile ilgili de olabilmektedir. Hastanın hastanede kaldığı süre boyunca (YBÜ dahil), aldığı ilaçların tutulduğu tablo ise PRESCRIPTIONS tablosudur. Bu tabloda, hastanın hangi ilacı aldı ğı, ilacın tipi, çeşitli tanımlayıcı numaraları, ilacına



Tarih	Olay
01/03/2017	Hastaneye Kabul
	Kan Ölçümü
	Tansiyon Ölçümü
02/03/2017	Tansiyon Ölçümü
	Kan Ölçümü
	Emar Çekimi
03/03/2017	Hastaneden Taburcu Olma
22/05/2017	Hastaneye Kabul
	Muayene
	Hastaneden Taburcu olma

Şekil 3.3: Örnek hasta ziyaret ve kayıt defteri.

zaman, hangi dozda aldığı gibi bilgiler yer almaktadır.

MIMIC-III içerisinde yer alan ve bu tez kapsamında kullanılan/kullanılmayan bütün tabloların isimleri, satır/sütun sayıları, ve kısa açıklamaları Çizelge 3.1’de okuyucu ile paylaşılmıştır.

3.2 MIMIC-III Ön İşleme Adımları

Yapılan deneylerde, temel olarak MIMIC-III veri seti içerisindeki gözlem verileri (yaşamsal veriler, lab sonuçları vb.), klinik notlar ve hastalara verilen ilaç (tedavi) bilgileri kullanılmaktadır. Bu verileri doğrudan kullanmak, ESK verisinin doğasına göre oldukça zordur. Sağlık alanında yapay öğrenme çalışmaları gerçekleştiren araştırmacıların karşılaştığı önemli problemlerden birisi, açık kaynak veri setleri için standart hale gelmiş veri ön işleme adımlarının bulunmamasıdır. ESK verisinin ham halinin çok farklı veri türlerini birden içermesi, bu verinin yapay öğrenme modellerine girdi olarak verilmesinden önce birçok veri ön işleme adımından geçirilmesine sebep olmaktadır. Yayınlanan çalışmaların büyük çoğunluğunda, çalışmaya ait kodların paylaşılmaması, veri ön işleme adımlarının tekrarlanması ile sonuçlanarak araştırmacılara ek maliyet getirmektedir. Veri ön işleme adımlarının standart hale getirilmemesi ve açık kaynak olmamasının bir diğer dezavantajı ise, yapılan çalışmaların tekrarlanmasını güçleştirmekle beraber, çalışma sonuçlarının doğrudan birbirleri ile karşılaştırılmasında zorlaştırmasıdır. Literatürdeki bu problemi çözmek adına MIMIC-III veri seti için çeşitli çalışmalar gerçekleştirilmiştir [84, 85, 86, 87]. Tez kapsamında yapılan deneylerde, MIMIC-III veri setinin kullanılabilmesi için veri ön işleme ve veriyi yapay öğrenme

Çizelge 3.1: MIMIC-III tabloları ve kısa açıklamaları.

Tablo İsmi	Satır Sayısı	Sütun Sayısı	Açıklama
ADMISSIONS	58,976	19	Hastaların, tekil hastane yatis, bilgisini içerir.
CALLOUT	34,499	24	YBÜ'nün taburcu planı ile ilgili bilgileri tutar.
CAREGIVERS	7,567	4	Tıbbi görevlilerin bilgileri içerir.
CHARTEVENTS	330,712,483	15	Yoğun bakımda kalan hastalara ait bütün tıbbi gözlem verilerini içerir.
CPTEVENTS	573,146	12	Hastalar üzerinde gerçekleştirilen prosedürlerin faturalandırılması için kullanılan Current Procedural Terminology (CPT) kodlarını içerir.
D_CPT	134	9	CPT kodları için genel tanımlarını içeren sözlük.
D_ICD_DIAGNOSES	14,710	4	International Classification of Disease (ICD-9) teşhis kodlarının bilgisini içeren sözlük.
D_ICD_PROCEDURES	3,898	4	Hastalara uygulanan tıbbi işlemler ile ilgili ICD-9 kodlarının bilgisini içeren sözlük.
D_ITEMS	12,487	10	YBÜ veritabanındaki "ITEMID" alanına ait bilgileri içeren sözlük (Laboratuvar ölçümleri hariç).
D_LABITEMS	753	6	Tüm laboratuvar ölçümleri için "ITEMID" alanına ait bilgileri içeren sözlük.
DATETIMEEVENTS	4,485,937	14	YBÜ içerisinde hastalara uygulanan tıbbi işlemlerin tarih kayıtlarını içerir.
DIAGNOSES_ICD	651,047	5	Hastalara konulan teşhis bilgisini (ICD kodları) içerir.
DRGCODES	125,557	8	Faturalandırma ve diğer işlemler için kullanılan Diagnosis Related Groups (DRG) kodlarının bilgisini içerir.
ICUSTAYS	61,532	12	Hastaların, tekil YBÜ yatis, bilgisini içerir.
INPUTEVENTS_CV	17,527,935	22	Philips CareVue klinik veri toplama sistemi ile takip edilen YBÜ hastalarının verisini içerir.
INPUTEVENTS_MV	3,618,991	31	iMDSoft Metavision klinik veri toplama sistemi ile takip edilen YBÜ hastalarının verisini içerir.
LABEVENTS	27,854,055	9	Hastane ve klinikler yapılan laboratuvar ölçüm verilerini içerir.
MICROBIOLOGYEVENTS	631,726	16	Hastaların mikrobiyoloji sonuç verilerini içerir.
NOTEEVENTS	2,083,180	11	Tıbbi görevlilerin hastalar için yazmış, notları, radyoloji raporları ve benzeri klinik notları içerir.
OUTPUTEVENTS	4,349,218	13	YBÜ'da yatan hastalara ait çıktıları içerir.
PATIENTS	46,520	8	Hastalara ait demografik ve benzeri bilgileri içerir.
PRESCRIPTIONS	4,156,450	19	Hastalara verilen ilaç bilgisini (reçeteler) içerir.
PROCEDUREEVENTS_MV	258,066	25	iMDSoft Metavision klinik veri toplama sistemi ile takip edilen YBÜ hastalarına uygulanan prosedür bilgisini içerir.
PROCEDURES_ICD	240,095	5	Hastalara uygulanan prosedürel bilgileri (ICD kodları) içerir.
SERVICES	73,343	6	Hastanın kabul edildiği/aktarıldığı servis bilgisini içerir.
TRANSFERS	261,897	13	Hastaların hastanede kaldığı süre boyunca fiziksel konularını (hangi serviste kaldığı gün ve transfer bilgisini) içerir.

modellerine girdi olabilecek şekilde dönüş, türme işlemi gerekmektedir. Bu işlemi gerçekleştirebilmek adına şu an için en güncel çalışma olan MIMIC-Extract [87] çalışmasından yararlanılmıştır. MIMIC-Extract çalışması ile MIMIC-III veri seti ön işleme tabii

tutulmuş, tutularak yapay öğrenme modelleri için zaman-serisi öznitelikleri yaratılmıştır. Zaman-serisi verileri haricinde kullanılan diğer veri türleri (klinik notlar, medikal terimler, ilaç isimleri) ise bu tez kapsamında çıkartılmış, tır. Bu verilerin çıkartılması ve ön işleme adımlarının detayları ise Bölüm 5 ve Bölüm 7’de okuyucu ile paylaşılmış, tır.

MIMIC-Extract, MIMIC-III veri seti için geliştirilmiş, açık kaynak bir veri ön işleme ve veriyi hazır hale getirme çalışmasıdır. ESK verilerinin gürültülü doğası gereği, veri içerisindeki birimleri standardize etmek, aykırı değerleri tespit etmek gibi işlemlerin yapılması gerekmektedir. MIMIC-Extract bu işlemleri gerçekleştirmekle beraber, zaman serisi verilerini saatlik olarak gruplayarak daha gürbüz (robust) nitelikte öznitelikler çıkarmayı amaçlamaktadır. Üzerinde deneylerin kolay yapılabilmesi adına hastanın hastane içinde mortalite (in-hospital mortality), yoğun bakımda mortalite (in-ICU mortality), yoğun bakımda 3 günden fazla kalma (LOS > 3) ve yoğun bakımda 7 günden fazla kalma (LOS 7) problemlerine uygun veri seti oluşturulmuştur.

MIMIC-Extract, ilk aşama olarak çeşitli kriterlere uymayan hastaları eleyerek üzerinde çalışılacak probleme daha uygun bir veri seti hazırlamaktadır. Bu eleme kriterleri MIMIC-Extract çalışması için kullanılacak olan araştırmacı tarafından değiştirilebileceği gibi (kod değişiminden ziyade sadece parametre değeri değiştirilerek bu işlem gerçekleştirilebilir.) çalışmada kullanılan kriterler aşağıda paylaşılmıştır:

- Hastaların en az 15 veya 15 yaşından büyük olmaları gerekmektedir.
- Birden fazla kez hastaneyi ziyaret etmiş, olan hastaların sadece ilk ziyaretindeki veriler kullanılmış, tır.
- Yoğun bakım ünitesinde en az 12 saat, en fazla 10 gün kalmış olan hasta verileri kullanılmış, tır.

Hasta seçimi yukarıda paylaşılan 3 ana kriter üzerinden gerçekleştirilmiştir, tır. Daha sonra ise bu hastalara ait yaşamsal gözlem verileri ve laboratuvar sonuçları ön işleme tabii tutulmuştur. ESK veri setleri içerisinde yaşamsal gözlem verileri ve laboratuvar sonuçları farklı ölçüm değerleri ile saklanabilmektedir. Bu birimleri standart hale getirmek için örneğin, tüm hasta kiloları kilogram, boyları santimetre, vücut suları ise santigrat birimine dönüştürülmüştür. Bu ve bunun gibi birim dönüşümleri istendiği takdirde kullanıcı tarafından kolayca gerçekleştirilebilmektedir. Birim standartlaştırma işleminden sonra veri içerisindeki aykırı değerleri bulmak adına işlem yapılmaktadır. Bu işlem için, alan uzmanı olan klinik uzmanların bilgileri ışığında yaşamsal gözlem verileri ve laboratuvar sonuçlarının hangi alt ve üst limitten değerleri alabileceği belirlenmiş ve bu değerlerin dışında kalanlar aykırı değer olarak işaretlenerek eksik veri olarak değerlendirilmiştir. Birimleri standartlaştırma ve aykırı değerleri çıkartma

işlemlerinden sonra öznitelikler saatlik olarak gruplanarak birleştirilmiştir. Hastalara ait öznitelikler farklı zamanlarda ölçüldüğü için, bu özniteliklerin doğrudan kullanılması ortaya seyrek vektörler çıkmasına neden olmaktadır. Bu durumu engelleyebilmek adına, yapılan deneyler ve yansımalar bulgu verileri saatlik bazda gruplanarak daha yoğun vektörel temsiller ile çalışma fırsatı yaratılmıştır. Son olarak ise, öznitelikler içerisinde klinik olarak birbiri ile ilişkisi içinde olanlar birleştirilmiş ve böylece daha gürbüz öznitelik kümesi oluşmasına olanak sağlamıştır. Örneğin, MIMIC-III veri seti içerisinde "KalpHızı(HeartRate)" değişkeni CareVue sistemi tarafından 211 nolu ID ile veritabanında tutulurken 2008 yılı sonrasındaki kullanılan MetaVision tarafından ise 220045 nolu ID ile tutulmaktadır. Bu iki tekil özneliğin birleştirilmesi örnek bir semantik gruplama şeklidir. Diğer ön işleme adımlarında olduğu gibi bu semantik gruplama da konfigürasyon dosyasından yönetilmekte ve istendiği zaman gruplama biçimi araştırmacı tarafından güncellenebilmektedir. Bu anlatılan adımlar özetle Şekil 3.4'de okuyucu ile paylaşılmıştır.

Tez kapsamında hastanın zamana bağlı öznitelikleri olarak da isimlendirilen, hastaya ait laboratuvar ve yansımalar gözlem verileri 104 adet öznitelik ile temsil edilmiş ve bu özniteliklerin isimleri Çizelge 3.2'da paylaşılmıştır.

MIMIC-Extract veri ön işleme adımları ile oluşan veri seti ile farklı klinik problemlerin çözümleri üzerinde çalışmalar yapılmaktadır. Bu tezde, literatürdeki önemli ve temel iki klinik problem üzerinde çalışılmıştır. Bu problemler, mortalite ve yoğun bakımda kalma süresini tahmin etmektir. Mortalite ve YBÜ'de kalma problemleri kendi içlerinde de ikiye ayrılmıştır. Toplamda üzerinde çalışılan dört problemin tanımları aşağıda paylaşılmıştır:

1. Hastane içi mortalite: YBÜ'sine giriş yaptıktan sonra hastane içerisinde ölen hastaları tahmin etme
2. YBÜ içi mortalite: YBÜ'sine giriş yaptıktan sonra YBÜ içerisinde ölen hastaları tahmin etme
3. YBÜ kalma süresi > 3 gün: YBÜ'de 3'den fazla süre kalan hastaları tahmin etme
4. YBÜ kalma süresi > 7 gün: YBÜ'de 7'den fazla süre kalan hastaları tahmin etme

MIMIC-Extract veri ön işleme adımları sonucunda elde edilen derlem içerisinde, MIMIC-III veri seti içerisindeki 46,520 hastadan 34,472 adet hastanın verisi kalmıştır. Bu hastaların demografik bilgileri incelendiğinde 19,498 tanesinin erkek, 14,988 tanesinin kadın hasta olduğu görülmektedir. Hastaların yaş dağılımı incelendiğinde ortalama

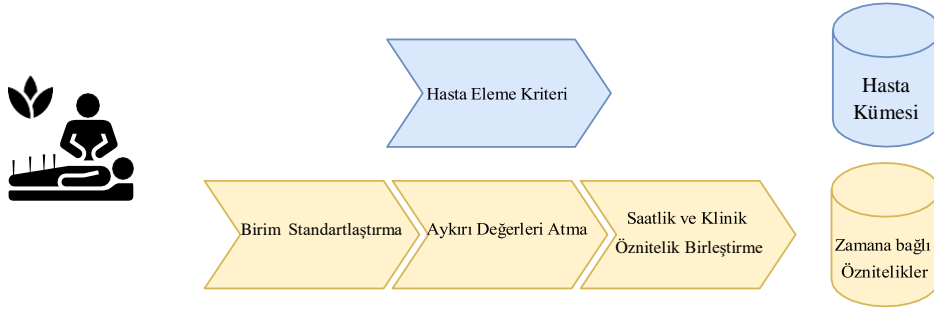
Çizelge 3.2: MIMIC-Extract çalışması ile çıkarılan ve deneylerde kullanılan 104 adet zaman-serisi öznitelik.

Öznitelik İsimleri		
alanine aminotransferase	fraction inspired oxygen set	plateau pressure
albumin	glasgow coma scale total	platelets
albumin ascites	glucose	positive and-expiratory pressure
albumin pleural	heart rate	positive end-expiratory pressure set
albumin urine	height	post void residual
alkaline phosphate	hematocrit	potassium
anion gap	hemoglobin	potassium serum
aspartate aminotransferase	lactate	prothrombin time inr
basophils	lactate dehydrogenase	prothrombin time pt
bicarbonate	lactate dehydrogenase pleural	pulmonary artery pressure mean
bilirubin	lactic acid	pulmonary artery pressure systolic
blood urea nitrogen	lymphocytes	pulmonary capillary wedge pressure
calcium	lymphocytes ascites	red blood cell count
calcium ionized	lymphocytes atypical	red blood cell count ascites
calcium urine	lymphocytes atypical csl	red blood cell count csf
cardiac index	lymphocytes bodyfluid	red blood cell count pleural
cardiac outputfick	lymphocytes percent	red blood cell count urine
cardiac output thermodilution	lymphocytes pleural	respiratory rate
central venous pressure	magnesium	respiratory rate set
chloride	mean blood pressure	sodium
chloride urine	mean corpuscular hemoglobin	systemic vascular resistance
cholesterol	mean corpuscular hemoglobin concentration	systolic blood pressure
cholesterol hdl	mean corpuscular volume	temperature
cholesterol ldl	monocytes	tidal volume observed
co2	monocytes csl	tidal volume set
co2 (etco2, pco2)	neutrophils	tidal volume spontaneous
creatinine	oxygen saturation	total protein
creatinine ascites	partial pressure of carbon dioxide	total protein urine
creatinine bodyfluid	partial pressure of oxygen	troponin-i
creatinine pleural	partial thromboplastin time	troponin-t
creatinine urine	peak inspiratory pressure	venous pvo2
diastolic blood pressure	ph	weight
eosinophils	ph urine	white blood cell count
fibronogen	phosphate	white blood cell count urine
fraction inspired oxygen	phosphorous	-

MIMIC-III

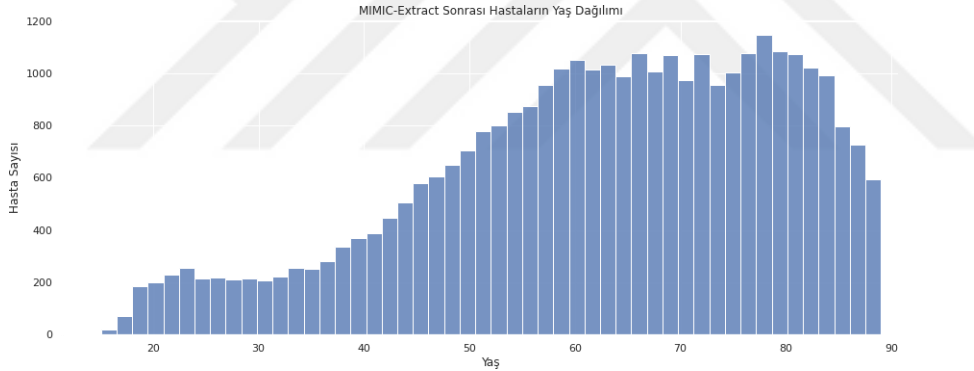
MIMIC-Extract

MIMIC-Extract Çıktısı



Şekil 3.4: MIMIC-Extract temel veri ön işleme adımları.

yaşın 62.34, standart sapmasının ise 16.91 olduğu görülmüştür. MIMIC-Extract veri ön işleme adımlarının yaratmış olduğu derlem içerisindeki hastaların yaş dağılımı Şekil 3.5’de gösterilmiştir. 34,472 hastanın ortalama yoğun bakımda kaldıkları gün sayısı 2.63, standart sapması 1.98 gün olmuştur. Hastaların genel olarak yoğun bakımda kaldıkları günlerin dağılımı Şekil 3.6’de okuyucular ile paylaşılmıştır.

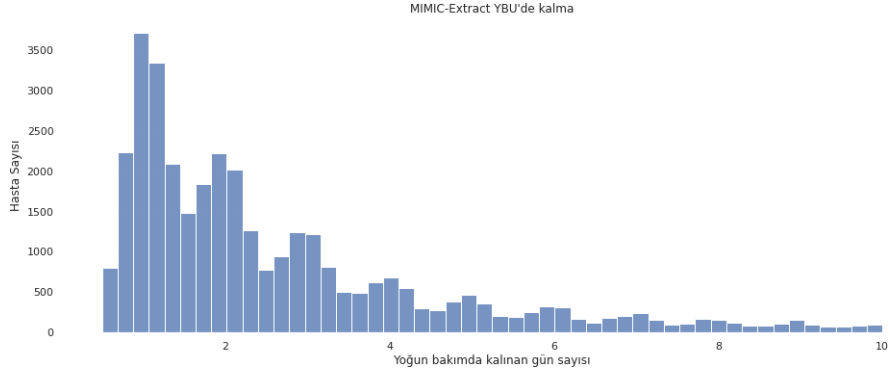


Şekil 3.5: MIMIC-Extract veri ön işleme sonrası, hastalarının yaş dağılımı.

MIMIC-III veri setinin orijinal durumundaki ve MIMIC-Extract veri ön işleme adımları uygulandıktan sonraki durumundaki hasta sayısı, hastane başvuru sayısı, ve yoğun bakım ünitesine başvuru sayıları Çizelge 3.3’de paylaşılmıştır. Tablo içerisindeki

Çizelge 3.3: MIMIC-III veri setinin, MIMIC-Extract uygulanmadan önceki ve sonraki veri istatistikleri.

	Hasta Sayısı	Hastane Başvuru Sayısı	YBÜ Başvuru Sayısı
MIMIC-III	46,520	58,976	61,532
MIMIC-III (>15 yaş, ından büyükler)	38,597	49,785	53,423
MIMIC-Extract	34,472	34,472	34,472
MIMIC-Extract (en az 24 saat YBÜ kalan hastalar)	23,937	23,937	23,937



Şekil 3.6: MIMIC-Extract veri ön işleme sonrası, hastalarının yoğun bakımda kalma sürelerinin dağılımı.

istatistiklerden görülebileceği üzere, MIMIC-Extract çalışmasının içerisinde yer alan kısıtlamalar (hasta yaşı sınırı, hastalara ait sadece ilk hastane ziyaret verilerinin ele alınması, ve diğer veri ön işleme adımları gibi) derlem içerisindeki hasta sayısının yarıyarıya azalmasına sebep olmuştur. Tez kapsamında yapılan çalışmalarda MIMIC-Extract'in 23,937 hasta sayısıyla temel alınarak çalışmalar gerçekleştirilmiştir. Buna rağmen çalışmalar özelinde de hasta eliminasyonu gerçekleştirilmiştir. Örneğin, çeşitli hastaların ilk 24 saatte klinik notlarının olmaması gibi bir durumda bu hasta derlem içerisinde o çalışmaya özel olarak çıkartılmıştır.



4. KULLANILAN YÖNTEMLER

4.1 ESK Verilerinin Analizinde Derin Öğrenme Yöntemleri

Elektronik sağlık kayıt verileri, içerisinde çok farklı veri türünü birlikte bulundurmaktadır. Yoğun bakımda yatan bir hastanın, zaman-serisi tabanlı yaşamsal gözlem verileri, çeşitli laboratuvar sonuçları verileri olabileceği gibi, yapısal olmayan klinik notları da bulunmaktadır. Farklı veri türlerini bir arada kullanarak basarılı bir şekilde temsil edebilmek ve bu temsiller üzerinden öznitelik vektörü çıkartabilmek adına yapılan deneylerde farklı yöntemler denenmiş, ve bu bölümde açıklanmıştır.

Sekans tabanlı (zamana bağımlı) verileri verimli bir şekilde işleyebilmek için, bu veri tipine özgü derin öğrenme tabanlı çeşitli algoritmalar önerilmiştir [69, 88]. Bu yöntemler sayesinde ise literatürde ses tanıma [89], müzik üretme [90], metin sınıflandırma [91], duygu sınıflandırma [92], makine çevirimi [93], varlık isim tanıma (NER) [94] gibi problemler için basarılı sonuçlar elde edilmiştir. Çalışmalar kapsamında kullanılan diğer veri türü olan metinsel veriler için ise öncelikle vektörel temsilleri öğrenilmeye çalışılmıştır. Bu kapsamda, Word2Vec [19], FastText [36], Doc2Vec [68], BERT [44], Sentence-BERT [95] yöntemlerinden yararlanılmış, ve bu yöntemlerin detayları okuyucuya aktarılmıştır. Bu temsil yöntemleri ile sayısal vektörler haline çevrilen metin verileri üzerinden öznitelik çıkarımı için ise 1D Evrimsel Sinir Ağları [38] ile gerçekleştirilmiştir. Son olarak ise, hastaya verilen ilaçların moleküler temsillerini elde edebilmek için kullanılan farklı yöntemlerin detayları (ECFP [58], MACCS [59], Mol2Vec [60], Smiles-Transformer [61]) bu bölüm içerisinde okuyucu ile paylaşılmıştır.

4.1.1 Tekrarlamalı yapay sinir ağı (RNN)

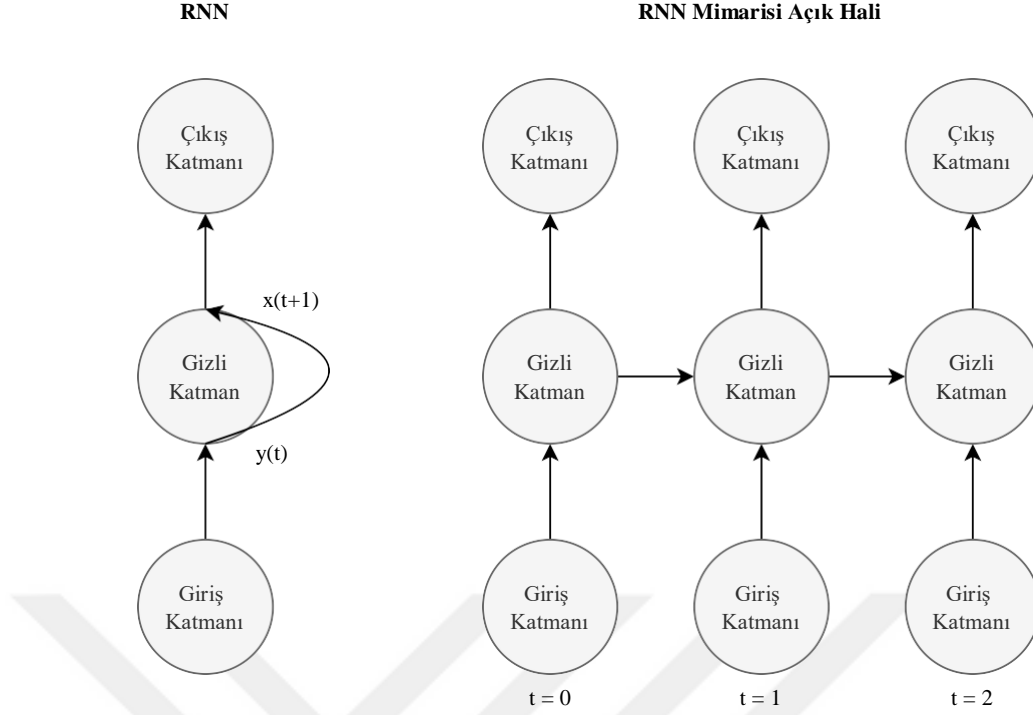
Zamana bağlı (sıra bilgisi içeren) verilerin klasik yapay sinir ağları ile işlenmesinin çeşitli zorlukları ve yetersizliklerini gidermek amacıyla Tekrarlamalı yapay sinir ağları (RNNs) [96] yöntemi önerilmiştir. Zamana bağlı verilerde, girdi ve çıktı uzunluklarının değişken olabilmesi, aynı zamanda zamanla öğrenilen bilginin aktarılabilmesi adına tekrarlamalı yapay sinir ağları bu iki durumu ele alabilecek şekilde tasarlanmıştır.

Zamana ve birbirine bağlı veriler üzerinde başarılı sonuçlar veren bu yöntem, metin verileri ile çalışılırken de kullanılmaktadır. Bunun temel sebebi, metinler içerisindeki

kelimelerin birbirinden bağımsız olmaması ve ilişki içerisinde olmasından kaynaklanmaktadır. RNN modelleri yapay sinir ağlarını kullanarak cümle içerisindeki geçmiş kelimeleri de hatırlayacak bir mimari sunmaktadır. Şekil 4.1’de RNN mimarisinin kapalı ve açık formu gösterilmiştir. Basitçe anlatmak gerekirse, gizli katman çıktısı sadece çıkış katmanına değil aynı zamanda yarattığı tekrarlanan döngü ile gizli katmana $t+1$ zaman için (diğer girdi ile beraber) aktarmaktadır. Bu yaklaşım, finans ve diğer başka alanlarda Auto-regressive moving average (ARMA) olarak da adlandırılmaktadır. RNN’lerin kullanmakta olduğu bu yaklaşım ilk bakışta karmaşık gelebileceği gibi temelde, $t=1$ anında girdi olarak gelen veriye ek olarak geçmişteki önemli bilgileri hatırlayabilmek adına $t=0$ anından gelen gizli katman verileri de girdi olarak verilmektedir. RNN modellerinin eğitilebilmesi için klasik yapay sinir ağları içinde geçerli olan geri yayılım (backpropagation) algoritmasının bir varyasyonu kullanılmaktadır. Klasik geri yayılım algoritmasının doğrudan kullanılamama sebebi RNN’ler içerisinde ek olarak zaman bilgisinin yer almasıdır. Bu sebeple zaman içerisinde geri yayılım (backpropagation through time - BPTT) yöntemi kullanılmaktadır. RNN modelleri içerisindeki yapıdan dolayı çok fazla öğrenilmesi gereken parametre olmamasına rağmen, özellikle uzun girdi türleri için eğitim süresinin maliyeti hızlı bir şekilde artmaktadır. Girdi uzunluğu arttıkça BPTT’nin güncellemesi için çok daha geçmişe giderek her zaman dilimi için yeniden türev hesaplanması gerekmektedir. Bu durum iki türlü probleme yol açmaktadır. Öncelikle hesaplama maliyeti artmakta ikinci olarak ise geriye doğru gradyan hesaplamalarında problemler yaşanmaktadır. Bu problemler gradyanların kaybolması (vanishing gradient problem) veya gradyanların patlaması (exploding gradient problem) olarak ikiye ayrılmaktadır. Hatanın geriye doğru hesaplanması sırasında türevlerin sürekli birbirleriyle çarpılması, geriye doğru gittikçe gradyanların çok küçülmesine veya çok büyümesine sebep olabilmektedir. Bu problemlerin önüne geçebilmek için literatürde Uzun Kısa Vadeli Hafıza Ağları (long short term memory networks, LSTMs) ve Geçitli Tekrarlayan Sinir Ağları (Gated Recurrent Unit, GRU) yöntemleri önerilmiştir. Tez kapsamında da doğrudan LSTM ve GRU yöntemleri kullanılmış, ve bu yöntemlerin detayları bir sonraki bölümde aktarılmıştır.

4.1.2 Uzun kısa vadeli hafıza ağları (LSTM)

Bir önceki bölümde anlatılan RNN yöntemi, girdi içerisindeki kısa süreli ilişkileri modelleyebilmesine rağmen, iki-üç zaman sonrasında gelen girdi ile geçmiş arasında bağlantı kurabilmekte zorlanmaktadır. Zaman serisi şeklinde bir girdi alındığında bu zaman içerisindeki bütün girdiyi hatırlayabilmek için RNN mimarisinin gelişmiş bir versiyonu olan Uzun Kısa Vadeli Hafıza Ağları (Long short term memory networks, LSTM) [69] Hochreiter ve Schmidhuber tarafından 1997 yılında önerilmiştir. LSTM



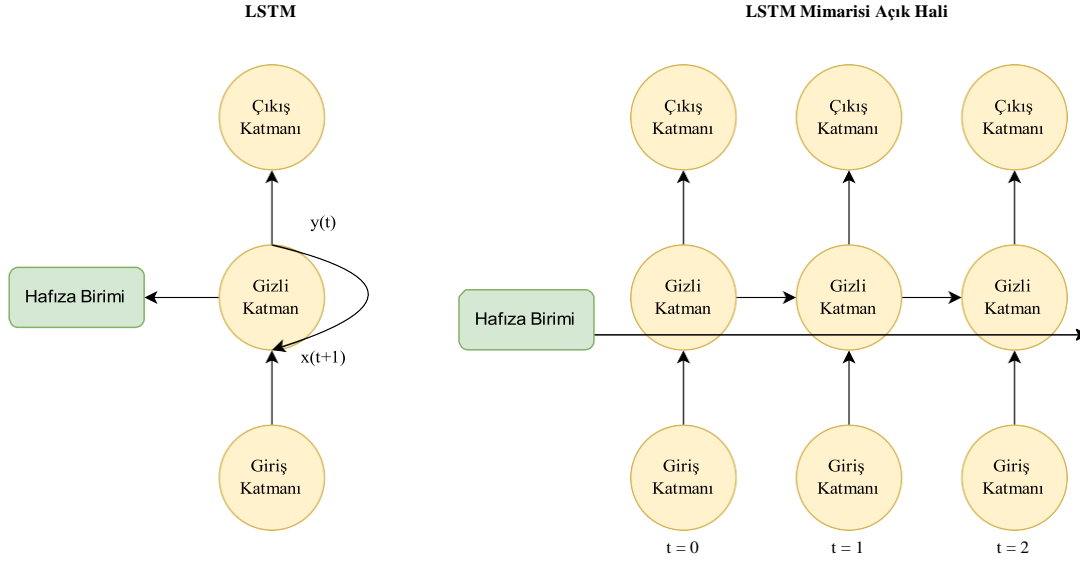
Şekil 4.1: Yinelemeli sinir ağı mimarisi.

yöntemi, her katmana durum (state) kavramı eklemeyi öne sürmüştür. Durum kavramı, hafıza (memory) görevi yaparak, modelin geçmişteki hangi bilgileri unutup hangi bilgileri hatırlaması gerektiği konusunda sorumludur. Bu sayede RNN yöntemlerinde karşılaşılan sadece bir iki zaman adım öncesini değil çok daha uzun zaman bağlantıları arasındaki ilişkileri yakalayabilmektedir. Şekil 4.2’de görüldüğü üzere, LSTM mimarisinin RNN mimarisinden temel farkı hafıza birimi içermesidir. Hafıza ünitesi içerisinde çeşitli kapılar ve matematiksel fonksiyonlar ile gizli hafıza biriminin unutacağı ve hatırlayacağı bilgiler güncellenmektedir. Bunun için unutma (forget), giriş (input) ve çıkış (output) kapıları kullanılır. Hücre yapısının detayları Şekil 4.3’de ve açıklamaları ise aşağıda okuyucu ile paylaşılmıştır.

Unutma Kapısı. Geçmiş, gizli katman bilgisi, h_{t-1} , ile yeni gelen girdi, x_t , sigmoid (σ) fonksiyonuna girdi olarak verilerek, bu bilgilerden hangilerinin unutulup hangilerinin unutulmayacağına karar verilmektedir. Bu fonksiyon içerisinde çıkacak 0 değeri, bilginin tamamen unutulacağı, 1 değeri ise bilginin tamamen saklanacağı anlamına taşımaktadır. Bu işlemi gerçekleştiren, unutma kapısının matematiksel formülü aşağıdaki gibidir.

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \quad (4.1)$$

Giriş Kapısı. Bu adımda ise yeni gelen bilginin hücrede saklanıp saklanmayacağına



Şekil 4.2: Uzun kısa vadeli hafıza ağı mimarisi.

karar verilir. Şekil 4.3'de i_t olarak sembolize edilen giriş kapısı, hangi değerlerin güncelleneceğine karar vermektedir. İkinci adımda ise, \tanh fonksiyonu yardımı ile \mathcal{C}_t ifadesi güncellenir ve hafızaya eklenebilecek yeni aday değerleri vektörü yaratılmaktadır.

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \quad (4.2)$$

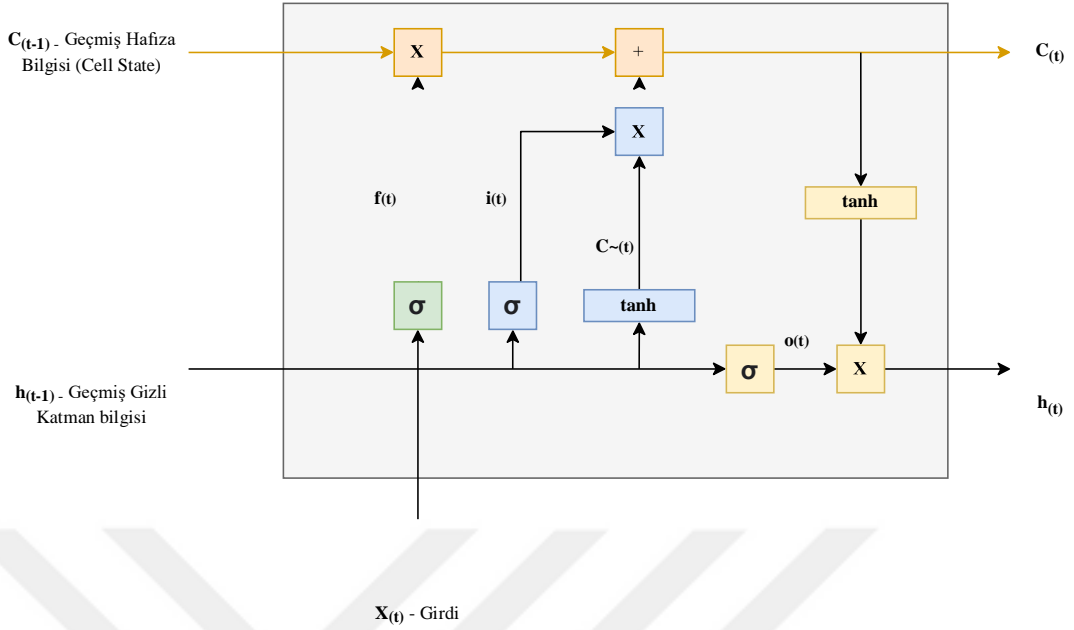
$$\mathcal{C}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c) \quad (4.3)$$

Bu adımlardan sonra, bir önceki zaman adımının C_{t-1} hafıza durumu güncellenerek, yeni hafıza durumu olan C_t ifadesinin güncellenmesi gerekmektedir. Bu işlem için, unutma kapısı ve giriş kapısı için hesaplanan güncel değerler kullanılmaktadır. Öncelikle geçmiş hafıza durumu C_{t-1} ifadesi ile unutma kapısında geri f_t çarpılmaktadır. Ardından yeni aday değer vektörü, \mathcal{C}_t ile bu değerlerin hangi miktarda güncellenmesi gerektiğine karar veren giriş kapısında geri olan i_t ifadesi çarpılır. Hesaplanan bu iki değer toplanarak güncel C_t değerine erişilebilmektedir. Anlatılan bu işlemin matematiksel karşılığı aşağıda ifade edilmiştir.

$$C_t = f_t \odot C_{t-1} + i_t \odot \mathcal{C}_t \quad (4.4)$$

Çıkış Kapısı. Son adımda ise, çıktılar ne üretilmesi gerektiğine karar vermek gerekmektedir. Bu sebeple, öncelikle hücre durumunun hangi bölümlerinin çıktısının alınacağına karar veren bir sigmoid katmanı çalıştırılmaktadır. Ardından, hafıza ünitesi \tanh aktivasyon fonksiyonundan geçirilmektedir. Bu iki katmanın çıktısı birbiri ile

LSTM Yapısı



Şekil 4.3: LSTM mimarisi hücre iç yapısı.

çarpılarak, sadece karar verilen kısımların çıktılar olarak alınması sağlanmaktadır.

$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o) \quad (4.5)$$

$$h_t = o_t \odot \tanh(C_t) \quad (4.6)$$

LSTM'de modelin eğitimi için kullanılan geriye yayılım algoritması RNN yönteminde kullanılan BPTT yöntemine benzemesine rağmen, LSTM içerisinde hafıza ünitesi sebebiyle, LSTM gradyanların kaybolması ve patlaması problemlerini yaşamamaktadır. LSTM içerisindeki hafıza ünitelerindeki ağırlık güncellemelerinin sadece önceki ve şu anki zaman değerleri ile hesaplanması, LSTM modelinin en azından gradyanların kaybolması probleminin önüne geçmesini sağlamaktadır. LSTM modelinin başarılı sonuçlar vermesiyle beraber, LSTM mimarisini geliştirme adına yeni varyantlar ortaya konulmuştur. Tez kapsamında zamana bağlı özellikleri klinik problemlerin tahmininde kullanmak için LSTM ile beraber, LSTM modelinin bir diğer varyantı olan GRU yöntemi kullanılmıştır. Bir sonraki bölümde GRU hakkında bilgi okuyucu ile paylaşılmıştır.

4.1.3 Geçitli tekrarlayan sinir ağı (GRU)

GRU yöntemi RNN tabanlı olup, LSTM mimarisinin bir başka varyantıdır. Ortaya çıkış sebebi, RNN modellerinin yaşamış olduğu gradyanların yok olması problemini

engellemeye ve LSTM'in ihtiyaç duyduğu hesaplama gücünü azaltmaya çalışmasıdır. 2014 yılında Cho vd. [88] tarafından önerilen bu yöntem, içerisinde LSTM benzeri kapılar bulundurmaktadır. LSTM mimarisinden temel farkı ise, içerisinde üç adet kapı (gate) yerine, güncelleme (update) ve sıfırlama (reset) isminde iki kapı bulundurmasıdır. Ayrıca GRU içerisinde, gizli katman bilgileri ve hafıza bilgileri ayrı ayrı tutulmak yerine birleştirilmiş, olarak bulunmaktadır. Bu mimari tasarım, GRU yöntemini başarımlı olarak hemen hemen LSTM ile aynı noktada tutarken, hesaplama hızı olarak LSTM yöntemine göre öne geçirmektedir. GRU modelinin yapısı Şekil 4.4'de okuyucu ile paylaşılmıştır. Güncelleme Kapısı. Girdi olarak, bir önceki zaman adımından aktarılan hafıza bilgisi olan h_{t-1} 'yi ve yeni gelen girdi bilgisi olan x_t 'yi kullanan güncelleme kapısı, ilgili girdileri kendi ağırlıkları olan W ve U matrisleri ile çarptıktan sonra sigmoid fonksiyonuna girdi olarak vererek, modelin geçmişteki bilgilerin ne kadarlık kısmını geleceğe taşımasını gerektiğini bulmaya çalışmaktadır. Bu işlemin matematiksel formülü aşağıda paylaşılmıştır:

$$z_t = \sigma(W_z \cdot x_t + U_z \cdot h_{t-1}) \quad (4.7)$$

Unutma Kapısı. Unutma kapısının ana görevi, geçmişteki bilginin ne kadarının unutulmasını gerektiğine karar vermektir. Güncelleme kapısı ile çok benzer bir işlem yapan unutma kapısında fark, çarpılan ağırlık matrisleridir.

$$r_t = \sigma(W_r \cdot x_t + U_r \cdot h_{t-1}) \quad (4.8)$$

Güncelleme ve unutma kapılarının değerlerinin hesaplanmasından sonra bu değerler, ilgili ağırlık matrisleri ile çarpılarak çıktıya katkısı hesaplanmaktadır. Bu sebeple, yeni hafıza birimi h_t^* , geçmiş bilgiler ve unutma kapısında geri ile beraber aşağıdaki gibi hesaplanmaktadır.

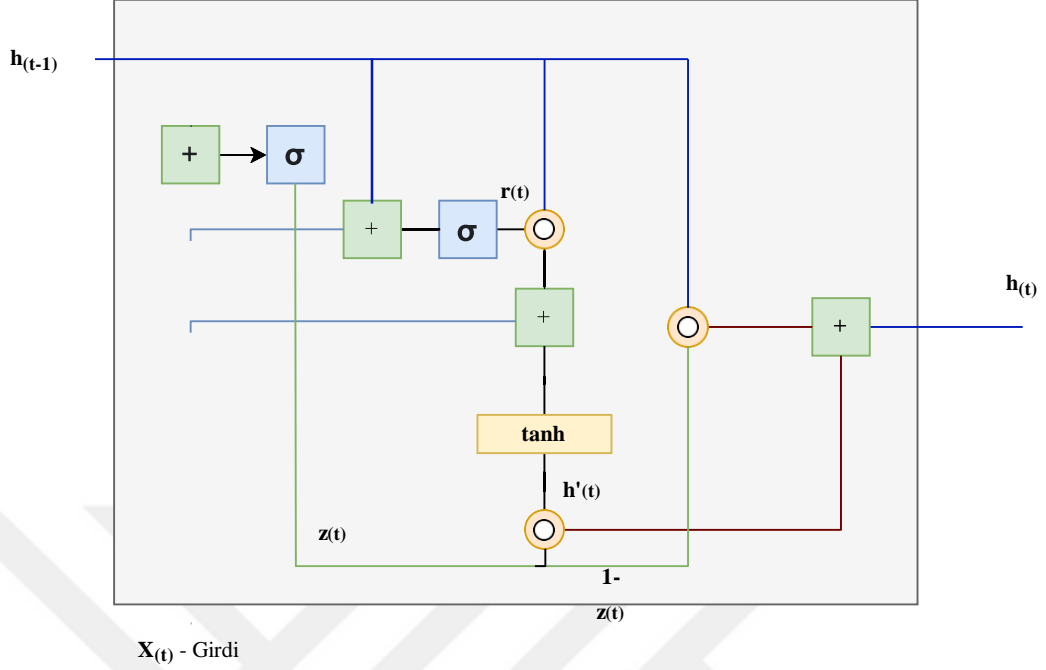
$$h_t^* = \tanh(W \cdot x_t + r_t \cdot U \cdot h_{t-1}) \quad (4.9)$$

Son olarak ise bir sonraki zaman damgasında kullanılacak olan hafıza bilgisi olan h_t 'yi hesaplamak için güncel hafıza bilgileri, güncelleme kapısının değeri ile çarpılarak hesaplanır.

$$h_t = z_t \cdot h_{t-1} + (1 - z_t) \cdot h_t^* \quad (4.10)$$

Özetle, GRU mimarisi LSTM mimarisine ile benzer bir amaçla RNN içerisindeki gradyanların yok olması ve uzun vadeli ilişkileri yakalayamaması problemlerini çözmek için önerilmiştir. LSTM ile benzer mantıkta çalışmasına rağmen temel bazı farklar mevcuttur. İlk olarak, LSTM, içerisinde giriş, çıkış, ve unutma isminde üç kapı bulundururken,

GRU Yapısı



Şekil 4.4: GRU hücre içi yapısı.

GRU mimarisi içerisinde güncelleme ve sıfırlama isminde iki kapı bulunmaktadır. İkinci olarak ise GRU, gizli katmanla taşınan ek bilgi haricinde LSTM mimarisindeki gibi ek bir hafıza ünitesi içermemektedir. Bu sebeplerden ötürü, GRU modellerinin genel olarak eğitim süreleri LSTM modellerinden daha hızlıdır. Tez kapsamında yapılan deneylerde hem LSTM hem de GRU modelleri denenmiş, model performans başarımları ve eğitim süreleri dikkate alınarak tercihler gerçekleştirilmiştir.

4.1.4 Evrimsel sinir ağları(CNN)

Evrimsel sinir ağları(CNN), insan beyninin görsel korteksinin çalışma prensibinden yola çıkarak geliştirilmiş, yapay sinir ağlarının varyantı olan bir öğrenme yöntemi olarak önerilmiştir. İnsan beyninin görme yetisinin nasıl çalıştığı ile ilgili çalışmalar 1960'lı yıllarda Hubel ve Wiesel tarafından [97, 98, 99] gerçekleştirildikten sonra bu çalışmalar yapay sinir ağlarına adapte edilmeye çalışılmıştır. Pratik uygulamaya dönüşen ilk örneklerden birisi, el yazısı ile yazılan rakamların tespit etme problemi üzerine olmuş, ve Lecun vd. tarafından 1989 yılında [100] gerçekleştirilmiştir.

2012 yılında Alex Krizhevsky tarafından önerilen AlexNet mimarisi [101] Evrimsel sinir ağları ve derin öğrenme alanının önemli bir kırılım noktası olmuştur. Bu mimaride ilk defa 8 katmanlı derin bir evrimsel sinir ağı önerilerek 62 milyon parametrelilik

bir model eğitilmiş ve ImageNet [102] yarışmasında bir önceki seneye göre %8'lik bir iyileşme göstermiştir. Daha sonra ise literatürde bu alanda oldukça fazla çalışma gerçekleştirilmiştir, VGGNet [103], InceptionNet [104], ResNet [105] gibi mimariler önerilmiş, ve hala aktif bir araştırma alanı olarak çalışılmaya devam etmektedir.

Evrışimsel sinir ağlarının görüntüler üzerine oldukça başarılı sonuçlar vermesi ile birlikte, bu yöntemin, zaman serisi [106], sinyal [107], veya metinsel veri [38] gibi girdi türleri üzerinde de etkisi araştırılmıştır. Bir boyutlu bu veriler üzerinde de başarılı sonuçlar veren 1D evrışimsel sinir ağları, bu tez kapsamında da klinik notlar ve ilaç temsilleri üzerinden öznetelik çıkartma işleminde kullanılmıştır.

Evrışimsel sinir ağları temelde iki farklı işlem gerçekleştirilmektedirler: evrişim (convolution) ve örnekleme (pooling). Genellikle CNN mimarileri, evrişim katmanı, örnekleme katmanı ve tam bağlantılı yapay sinir ağı ile devam eden mimariye sahiptir. Bu mimari içerisinde evrişim katmanı ve tam bağlantılı yapay sinir ağı içerisinde öğrenilmesi gereken parametreler bulunurken, örnekleme katmanında herhangi öğrenilecek bir parametre bulunmamaktadır. CNN içerisinde bulunan işlemler aşağıda detaylı bir şekilde açıklanmıştır.

Evrışim. Evrişim, bir görüntü üzerine filtre uygulayarak, görüntü içerisindeki kalıpları yakalamaya gılayan işlemdir. Şekil 4.5'de örnek bir görüntü matrisi ve filtre matrisi verilmiş olup, eleman bazında çarpma işlemi uygulanarak nokta çarpımı gerçekleştirilmiştir. Ortaya çıkan 2×2 matrisi ise, uygulanan filtrenin orijinal görüntü içerisinde yoğun olduğu kısımları göstermektedir. Bu işlem temel olarak CNN içerisindeki evrişim katmanında gerçekleştirilmektedir. Her bir evrişim katmanı içerisinde birden çok filtre, bir diğer bilinen ismi ile çekirdek (kernel) bulunmaktadır. Filtreler içerisindeki ağırlıklar ise modelin öğrenmesi gereken değerlerdir. Herhangi $n \times n$ boyutundaki bir görsele $f \times f$ boyutunda bir filtre uygulandı ğında aşağıdaki formüle göre çıktı boyutu hesaplanabilmektedir.

$$(n-f+1) \times (n-f+1) \quad (4.11)$$

Dolgulama (Padding) ve Adım Kaydırma (Stride). Görüntü ile filtre arasındaki çarpımlar esnasında boyut uyumsuzlu ğı sebebi ile görüntünün kenar kısımları ile filtre arasında çarpım işleminin gerçekleştirilemedi ği durumlar olabilmektedir. Bu sebeple dolgulama işlemi ile görüntünün kenarlarına sıfır eklenerek bu problemin önüne geçilmektedir. Bir boyutlu bir vektör için örnek işlem Şekil 4.6'de gösterilmiştir. Adım kaydırma yöntemi ise, uygulanan filtrenin, görüntü matrisi üzerinde kaç adım birden kayacağını belirleyen parametredir. Hem dolgulama hem de adım kaydırma parametrelerine göre oluşacak olan sonuç matrisinin boyutu aşağıdaki formüllere göre hesaplanabilmektedir.

Örnek Görüntü Matrisi

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

Filtre

1	0	-1
1	0	-1
1	0	-1

Sonuç

4	-9
-1	-3

$$5 * 1 + 2 * 0 + 4 * -1 + 7 * 1 + 5 * 0 + 6 * -1 + 4 * 1 + 3 * 0 + 2 * -1 = 4$$

$$2 * 1 + 4 * 0 + 9 * -1 + 5 * 1 + 6 * 0 + 3 * -1 + 3 * 1 + 2 * 0 + 7 * -1 = -9$$

$$7 * 1 + 5 * 0 + 6 * -1 + 4 * 1 + 3 * 0 + 2 * -1 + 1 * 1 + 3 * 0 + 5 * -1 = -1$$

$$5 * 1 + 6 * 0 + 3 * -1 + 3 * 1 + 2 * 0 + 7 * -1 + 3 * 1 + 5 * 0 + 4 * -1 = -3$$

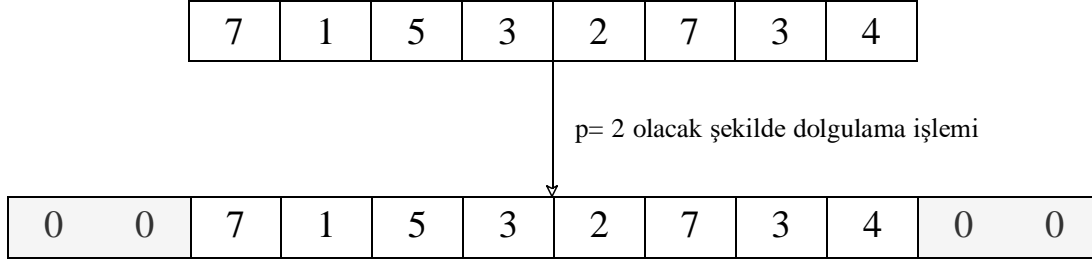
Şekil 4.5: CNN'de evrişim işlemi.

$$\text{Dolgu} = (n+2p-f+1) \times (n+2p-f+1) \quad (4.12)$$

$$\text{Adım kaydırma} = \left(\frac{n+2p-f}{s} + 1 \right) \times \left(\frac{n+2p-f}{s} + 1 \right) \quad (4.13)$$

Örnekleme (pooling). Evrişim operasyonunda filtreler kullanılarak veri içerisindeki öznitelikler çıkartılmaya çalışılmıştır. Önemli öznitelikleri elde etmeye çalışırken kullanılan filtre ve evrişim katmanlarının sayısının artması, hesaplanması gereken parametre sayısını artırarak, modelin eğitim süresini ve maliyetini artırmaktadır. Evrişimsel operasyon sonucunda oluşan matrisin boyutunu belirli bir noktada tutmak, eğitim süresini hızlandırmak ve öznitelik vektörlerinde sadece önemli kısımları elde etmek için örnekleme işlemi önerilmiştir. Örnekleme operasyonu, uygulandığı matris içerisinde bir alt örnekleme yaparak, yalnızca önemli ve ilgili bilgilerin korunmasını sağlamaktadır. Bu sayede, elenmesi istenen bilgiler elenerek gereksiz hesaplama maliyetinden kurtulmaktadır. Ayrıca bu yöntemin kullanılması modelin aşırı öğrenmesini de engellemeye yardımcıdır. Maksimum dolgu, ortalama dolgu gibi farklı dolgu teknikleri Şekil 4.7'de gösterilmiştir.

Evrişim işlemi, yukarıda anlatıldığı üzere, görüntülerdeki çeşitli ilişkileri yakalayabilmektedir. Görüntüler ile aynı olmamasına rağmen, metinler içerisinde de sıralı bir ilişki mevcuttur. Kelimeler arasındaki ilişkilerin farkı ise, görüntüdeki gibi 2 boyutlu olmamasıdır. Kelimeler arası ilişki üst/alt satır arasından ziyade tek boyutta yani sağ/sol komşular arasında mevcuttur. Bu nedenle, metin verileri için tek boyutlu uzamsal ilişkiler 1D evrişimsel katman ile gerçekleştirilebilmektedir. Şekil 4.8'de örnek bir 1D evrişim ve maksimum örnekleme işlemi gösterilmiştir. Her kelimenin 3 boyutlu



Şekil 4.6: CNN’de dolgulama (padding) işlemi.

vektörler ile temsil edildiği ve toplamda 6 kelimedenden oluşan bir cümleye 1D evrişim operasyonu uygulandığında, 6 boyutlu bir öznitelik vektörü elde edilmektedir. Şeklin sadeliğini koruyabilmek adına tek bir filtre ile işlem yapıldığı gösterilmiştir. Duygu analizi, soru sınıflandırma gibi metin tabanlı problemlerde 1D evrişimsel sinir ağlarının başarılı sonuçlar verdiği, Kim [38] tarafından 2014 yılında gösterilmiştir. Bu çalışma kapsamında, medikal terimlerin temsilleri üzerinden öznitelik çıkarım problemini çözmek için detayları Bölüm 5’de anlatıldı. Güzere, içerisinden birden çok evrişim katmanının olduğu 1D evrişimsel sinir ağları kullanılmıştır.

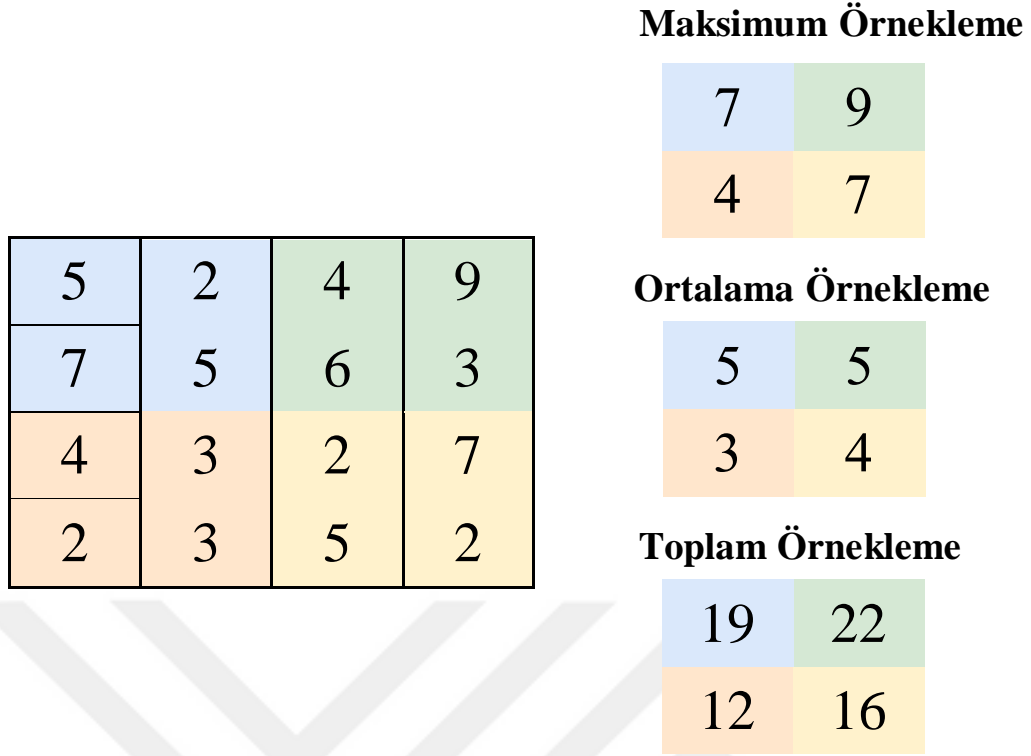
4.2 Kelime Temsil Yöntemleri

Bölüm 5 ve Bölüm 6’da gerçekleştirilen deneyler içerisinde hasta için yazılan klinik rapor ve notlar, eğitilen modellere ek bir veri türü olarak verilmiş ve model

başarımlarının artırılması hedeflenmiştir. Metinsel ifadelerden oluşan klinik notların yapay/derin öğrenme algoritmalarına girdi olarak verilebilmesi için öncelikli olarak vektörel ifadelerle dönüştürülmeleri gerekmektedir. Literatürde uzun yıllardır çalışılan, metinsel verilerin vektörel karşılıklarının öğrenilmesi konusu üzerine birçok yöntem geliştirilmiştir. Tez kapsamında kullanılan metin temsil yöntemleri ise Word2Vec [19], FastText [36], Doc2Vec [68], BERT [44], ClinicalBERT [46] ve Sentence-BERT [95] olmuştur. Bu bölümde bu yöntemlerin detayları aktarılmaktadır.

4.2.1 Word2Vec

Anlamsal olarak benzer kelimelerin temsil vektörlerini birbirine yakın tutarken, birbirine anlamsal olarak benzemeyen kelimelerin temsillerini birbirlerinden uzak tutmayı amaçlayan kelime temsili öğrenme algoritmalarından biri olan Word2Vec yöntemi, 2013 yılında Mikolov vd. [19] tarafından önerilmiştir. Word2Vec yönteminin temeli, kelimelerin derlem içerisinde beraber geçme bilgisi üzerine kurulmuştur. Belirli bir pencere içerisinde beraber geçen kelimelerin anlamsal olarak yakın olma olasılığı, kelimelerin matematiksel ifadeleri (vektörel) arasında da yakınlık kurulabilmesini sağlamaktadır. Önerilen bu yöntemde kullanılan yapay sinir ağı, giriş katmanı, gizli

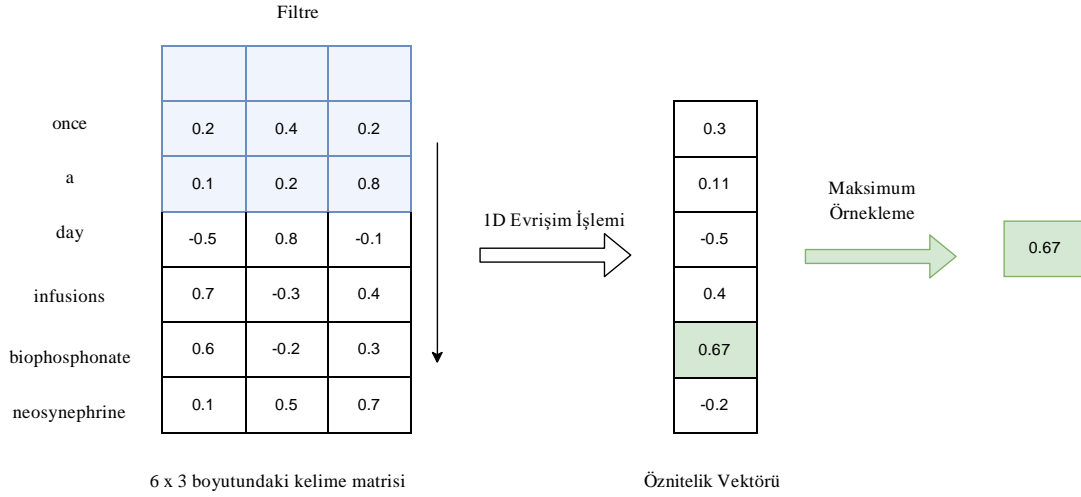


Şekil 4.7: CNN’de örnekleme (pooling) işlemi.

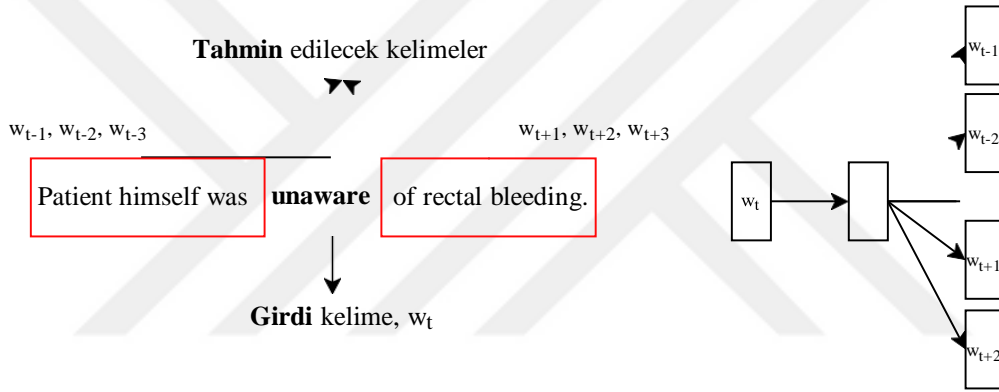
katman ve çıkış katmanından oluşan sığ (shallow) bir yapay sinir ağıdır. Word2Vec kelimelerin vektörel temsillerini öğrenirken, CBOW (Continuous Bag of Words) ve Skip-Gram isminde iki farklı yöntem kullanabilmektedir. CBOW yönteminde, bir kelimenin belirli bir penceredeki komsu kelimelerini kullanarak o kelime tahmin edilmeye çalışılmaktadır. Skip-gram yönteminde ise tam tersi olarak seçilen bir kelimedenden, komsu kelimelerin tahmin edilmesi gerçekleştirilmektedir. Bu iki yöntemin detaylarına şurada anlatılmaktadır.

Skip-Gram. Bu yöntem Şekil 4.9’de gösterildiği gibi, model girdi olarak bir kelimeyi alırken, girdi olarak verilen kelimenin komsu kelimelerini tahmin etmeye çalışmaktadır. Girdi, ve gizli katman olmak üzere iki katmandan oluşan yapay sinir ağının gizli katmanında n adet nöron bulunduran bu yapıdaki n ifadesi öğrenilecek olan kelime vektör boyutunu göstermektedir. Girdi ve çıkış katmanında ise m adet nöron bulunmakta ve bu m ifadesi derlem içerisindeki tekil kelime sayısını belirtmektedir. Önerilen modelin son katmanında ise aktivasyon fonksiyonu olarak, üretilen sonucu 0-1 arasına sıkıştırılan ve bütün çıktılarının toplamını 1 olacak şekilde formüle eden softmax kullanılmıştır.

CBOW. Bu yöntem ise skip-gram yönteminin aksine, komsu kelimeleri kullanarak, merkezdeki kelimeyi tahmin etmeye çalışmaktadır. Şekil 4.10’de örnek bir cümle ve model mimarisi gösterilmiştir. CBOW ile Skip-gram arasında yöntemsel fark olması



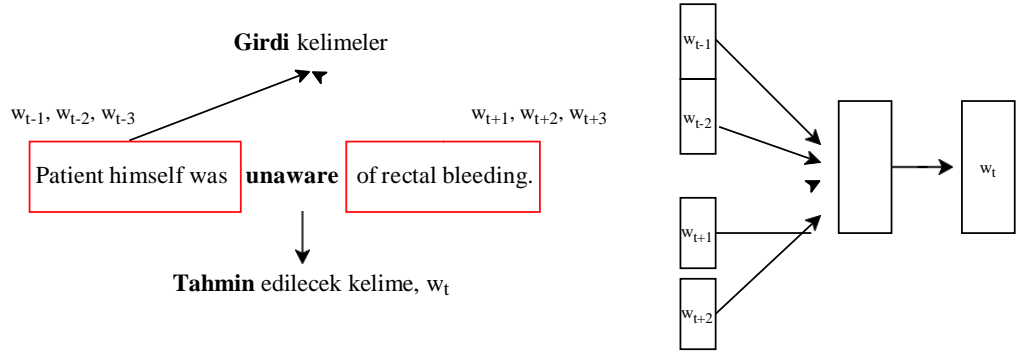
Şekil 4.8: Metinsel veri üzerinde 1-boyutlu evrişim ve maksimum örneklem işlemleri.



Şekil 4.9: Skip-Gram yönteminin gösterimi.

sebebiyle birbirlerine karşı avantajlı oldukları için, ilgili noktalar bulunmaktadır. Word2Vec modelini öneren Mikolov vd. da çalışmalarında bahsettikleri üzere, skip-gram yöntemi küçük derlemlerde, CBOW yöntemine göre daha başarılı çalışmaktadır. CBOW yöntemi ise özellikle sık geçen kelimeler için daha başarılı temsiller üretmekle beraber eğitim süresinde de skip-gram yöntemine göre daha hızlı çalışmaktadır.

Alana özel problemlerin çözümü için metin verileri kullanılmak istendiğinde, bu metinlerin alana özel eğitilmiş modeller ile vektörize edilmesi gerekmektedir. Google'un ve diğer araştırmacıların büyük derlemler ile eğittikleri genel Word2Vec modelinde, üzerinde çalışılan alana özel kelimelerin bu derlem içerisinde geçmeme veya çok az geçme ihtimalinden dolayı, alana özgü kelimelerin temsillerinde başarımların düşüklüğü yaşanabilmektedir. Tez kapsamında yapılan çalışmalar MIMIC-III veri seti üzerinde gerçekleştirildiği için, kullanılacak olan Word2Vec modelinin de bu veri seti üzerinden (en azından klinik alan verileri) eğitilmiş olması önemlidir. Bu kapsamda, Huang vd. [48] tarafından MIMIC-III içerisindeki klinik notlar kullanılarak eğitilmiş Word2Vec modeli



Şekil 4.10: COW yönteminin gösterimi.

kullanılmıştır. Bu model eğitiminde MIMIC-III klinik notları içerisindeki 2.8 milyar kelime kullanılmış ve her kelime için 100 boyutlu kelime temsilleri öğrenilmiştir. Modelin temsili gösterimi Şekil 4.11’de paylaşılmıştır.

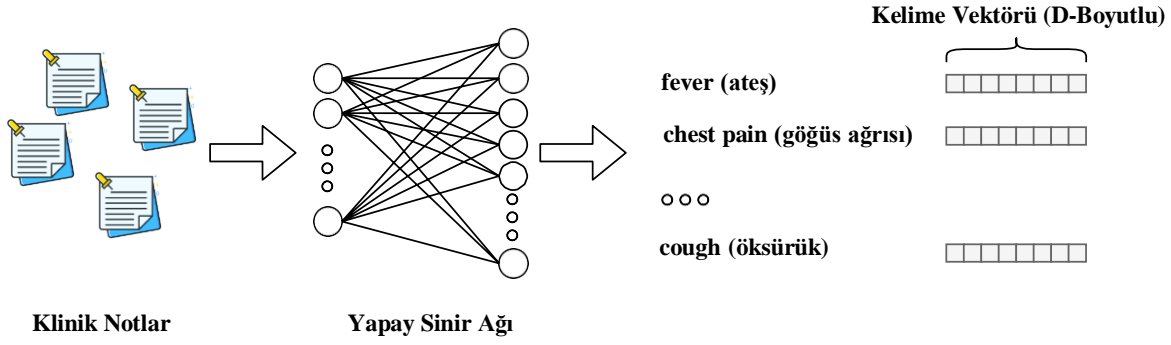
4.2.2 FastText

2013 yılında geliştirilen Word2Vec yöntemi ile doğal dil işleme alanında birçok farklı çalışma gerçekleştirilmeye başlanmıştır. Buna çalışmalarda yaşanan en önemli problemlerinden bir tanesi, Word2Vec yönteminin eğitildiği derlem içerisinde yer almayan bir kelimenin temsili çıkarılamaması olmuştur. Word2Vec yönteminin gelişmiş bir versiyonu olan FastText [36] yöntemi 2016 yılında Facebook içerisinde çalışan araştırmacılar tarafından önerilmiştir. FastText yöntemi, Word2Vec yönteminde olduğu gibi komsu kelimeleri tahmin etmek yerine komsu n-karakterlik ifadeleri tahmin etmeye çalışmaktadır. Örneğin, "clinical" kelimesi 2 ve 3 karakterlik ifadelerle bölünerek "cl", "cli", "li", "lin", "in", "ini", "ni", "nic", "ic", "ica", "ca", "cal" ifadelerine dönüşmektedir. Bu sayede, kelimelerin iç düzeni hakkında bilgi edinilebilmekte, Word2Vec yönteminin dezavantajlı olduğu ve daha önce görmediği kelimelere cevap verememe problemi ortadan kaldırabilmektedir. Word2Vec modelinde olduğu gibi tez kapsamında kullanılan FastText modeli, [48] çalışması tarafından MIMIC-III içerisindeki klinik notlar kullanılarak eğitilmiştir. Deneyler içerisinde kullanılan daha önceden eğitilmiş olan Word2Vec ve FastText modeline ilgili linkten ² ulaşılabilir.

4.2.3 Doc2Vec

Kelime temsillerinin öğrenilmesine yönelik önceki bölümlerde anlatılan Word2Vec ve FastText yöntemlerine ek olarak doğrudan doküman ve cümle temsillerini öğrenebilmek için Doc2Vec [37] yöntemi de tez kapsamında gerçekleştirilen deneylerde kullanılmıştır. Doc2Vec yönteminin isminden de anlaşılacağı üzere, yöntem, doküman

²<https://github.com/kexinhuang12345/clinicalBERT>

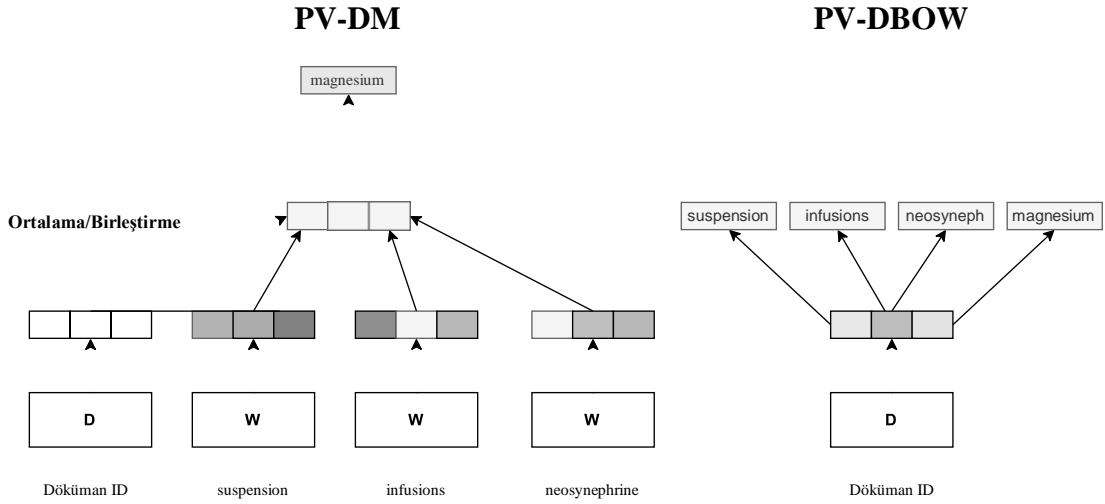


Şekil 4.11: Klinik alana özgü Word2Vec eğitimi.

veya paragraf seviyesinde bağlamsal vektörel temsil öğrenilmesini sağlamaktadır. Bir önceki bölümde anlatılan Word2Vec yöntemi kelime seviyesinde temsilleri öğrenmesine rağmen, doğal dil işleme alanındaki problemlerde genellikle birden çok kelimenin bir arada kullanıldığı metinler ile çalışmalar gerçekleştirilmektedir. Bu sebeple, doküman temsili yapılırken, Word2Vec ile temsil edilen kelimelerin ortalamalarını kullanmak yerine Le ve Mikalov tarafından Word2Vec yöntemini genişleterek, farklı uzunlukta dokümanlar için sabit uzunlukta vektörel temsil öğrenen Doc2Vec önerilmiştir. Doc2Vec yöntemi içerisinde, Word2Vec yönteminde olduğu gibi iki farklı öğrenme mimarileri mevcuttur. Bunlar; Distributed Memory Model of Paragraph Vectors (PV-DM) ve Distributed Bag-of-Words Models of Paragraph Vectors (PV-DBOW) yöntemleridir. PV-DM yöntemi, önceki bölümde detayları anlatılan Word2Vec içerisindeki CBOW yöntemi ile benzerlik göstermektedir. Bu yöntemde, doküman vektörleri, kelime vektörleri ile birleştirilerek veya ortalaması alınarak hedef kelime tahmini yapılmaktadır. PV-DBOW yöntemi ise Word2Vec yöntemindeki skip-grap yaklaşımına benzemektedir. Kelime temsilleri dikkate alınmadan, paragraf vektörünün, paragraf içerisinde rastgele seçilen kelimeleri tahmin etmesi beklenmektedir. Bu yöntemin eğitimi sırasında, gradyan inişi (gradient descent) ve geri yayılım (backpropagation) yöntemleri ile paragraf vektörleri öğrenilmektedir. Her iki yöntemin temsili, Şekil 4.12’de gösterilmiştir.

4.2.4 Bidirectional Encoder Representations from Transformers (BERT)

Yapay sinir ağları, içerisinde hafıza birimi bulunan özyinelemeli sinir ağları (Recurrent Neural Networks, RNNs), evrimsel sinir ağları (Convolutional Neural Networks, CNNs), dikkat mekanizmalı ağlar (attention mechanism) gibi yöntemler literatürdeki birçok problem üzerinde oldukça başarılı sonuçlar vermişlerdir. Bu yöntemlerin hesaplama maliyetinin fazlalığı, RNN tabanlı yöntemlerin paralelleştirilmesinde zorluklar gibi sebeplerden ötürü, Aralık 2017’de Vaswani vd. [45] tarafından Transformer tabanlı, tamamen yeni bir yaklaşım öneren bir yöntem geliştirilmiştir. Bu yöntem sayesinde

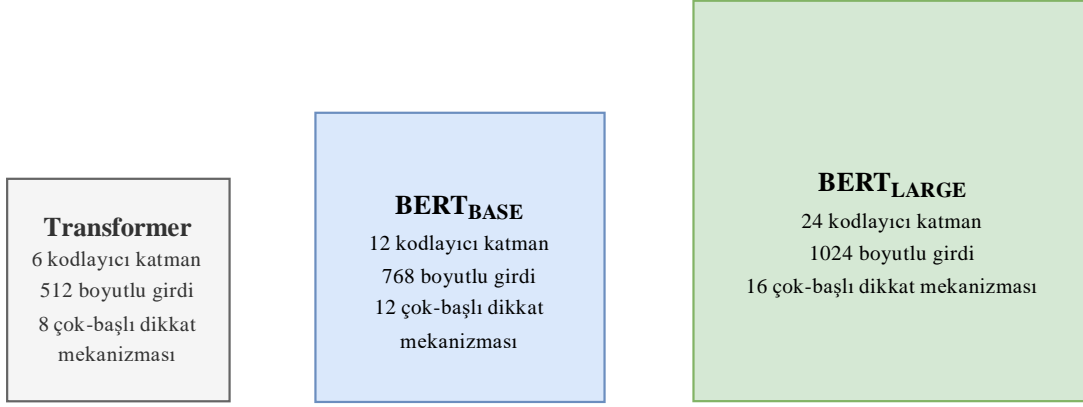


Şekil 4.12: Doc2Vec yöntemi içerisindeki farklı öğrenme mimarileri olan PV-DM ve PV-DBOW metodlarının gösterimi.

literatürde birçok yapay zeka probleminde en iyi sonuçlar alınmaya başlanmıştır. Transformer mimarisi, içerisinde birçok farklı yapıyı bir arada bulundurmaktadır ve bu yapılar kolayca değiştirilebilir olarak tasarlanmıştır. Bu yapılar; kodlayıcı (encoder), çözücü (decoder), gömme katmanı (embedding layer), pozisyonel kodlayıcı (positional encoding), çok-başlı dikkat mekanizması (multi-head attention), maskelenmiş, çok-başlı dikkat mekanizması (masked multi-head attention), toplayıcı ve normalize edici (add & norm), ileri beslemeli yapay sinir ağıdır.

Bidirectional Encoder Representations from Transformers (BERT) [44] çalışması ise 2018 yılında Devlin vd. tarafından sunulmuştur. Doğal Dil İşleme (NLP) alanında devrim niteliğinde olan, öğrenim aktarımı (transfer learning) imkanı sunan bu çalışma ile beraber NLP alanındaki hemen her problemde şimdiye kadar alınmış, en iyi sonuçlar alınmıştır. BERT, 2017 yılında çıkan transformer tabanlı mimariye, iki-yönlü çok-başlı dikkat mekanizması (bidirectional multi-head attention) ekleyerek modifiye etmiştir. Yeni eklenen bu yöntem ile, model bir cümleyi işlemeye çalıştığında sadece soldan sağa değil, aynı zamanda sağdan sola da işleyerek, cümlenin tamamına bakabilme yeteneği kazanmıştır.

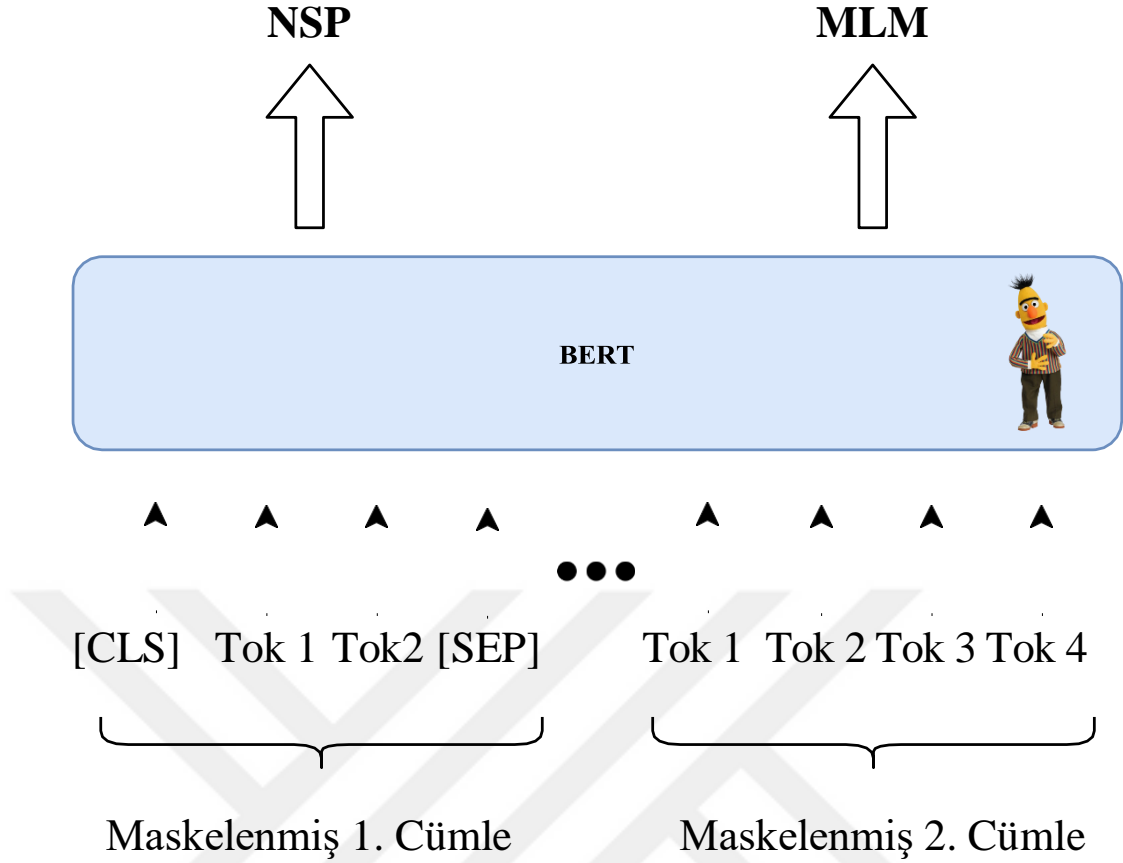
BERT modeli, transformer mimarisinden başka noktalarda da ayrılmaktadır. İlk olarak BERT, sadece kodlayıcı (encoder) katmanını içermektedir. Ayrıca bu kodlayıcı katmanları orijinal transformer mimarisi içerisinde bulunan kodlayıcı katmanlarına göre daha büyük bir yapıdadır. İki farklı mimari ile ortaya çıkan BERT, $BERT_{BASE}$ ve $BERT_{LARGE}$ olarak ikiye ayrılmaktadır. $BERT_{BASE}$, 12 kodlayıcı katmanının birleşmesi ve 12 çok-başlı dikkat mekanizmasından oluşurken, $BERT_{LARGE}$ ise 24 kodlayıcı katman ve 16 çok-başlı dikkat mekanizmasından oluşmaktadır. Bu bilgiler özet halinde Şekil 4.13'de



Şekil 4.13: Transformer, $BERT_{BASE}$ ve $BERT_{LARGE}$ yöntemlerinin mimari büyüklük karşılaştırılması.

okuyucu ile paylaşılmıştır.

BERT modeli kullanımı için iki aşamalı bir yapı önerilmiştir. Bunun için, Şekil 4.14’de gösterildiği üzere modelin ilk olarak ön eğitimden (pre-training) geçirilmesi daha sonra ise üzerinde çalışılmak istenen problem ile ince ayar (fine-tuning) yapılması gerekmektedir. BERT modelini ön eğitimden geçirmek için ise iki yöntem önerilmiştir, Maskeli Dil Modelleme (Masked Language Modeling, MLM) ve sonraki cümleyi tahmin etme (Next Sentence Prediction, NSP). MLM yönteminde, bütün tokenlerin %15’lik kısmı rasgele olarak seçilerek maskelenmektedir. Bu yöntem basit ve mantıklı olmasına rağmen, ince ayar yapılacağı zaman, [MASK] ifadesinin olmaması problem yaratabilmektedir. Bunun önüne geçebilmek adına, rastgele seçilen %15’lik kelimelerden, %80’i [MASK] ifadesi ile ($15\% * 80\% = 12\%$), %10’u rastgele bir kelime ile ($10\% * 15\% = 1.5\%$) ve %10’u da değiştirilmeden ($10\% * 15\% = 1.5\%$) bırakılmaktadır. MLM görevinde, modele birden çok cümle yerine, tekil cümleler verilerek eğitim gerçekleştirilmektedir. Soru-cevaplama (question-answering) gibi problemlerde ise BERT modeline birden çok cümle verilerek aralarındaki ilişkinin yakalanması beklenmektedir. Bunu sağlayabilmek için NSP yöntemi geliştirilmiştir. Modele girdi olarak A ve B cümleleri beraber verilirken, bu ikili cümlelerin %50’sinde ikinci cümlenin ilk cümleden bir sonraki cümle olması (isNext etiketi atanır), diğer %50’sinde ise olmamasına (NotNext etiketi atanır) göre girdi düzeni ayarlanır. Daha sonra ise, yeni bir ikili cümle verildiğinde, ikinci cümlenin ilk cümleden sonra gelip gelemeyeceğini model tahmin etmektedir. Ayrıca, NSP yönteminde BERT modeline girdi olarak iki cümle, iki yeni token ile birlikte verilmektedir. [CLS] tokeni ikili sınıflandırma tokeni olup, ilk sekansın başına eklenirken, [SEP] tokeni, sekansın bittiğini ifade etmektedir.



Şekil 4.14: BERT modelinin ön eğitim aşamasının gösterimi.

4.2.5 Clinical BERT

Bağlamsal kelime temsil yöntemlerinden BERT, bir önceki bölümde anlatıldığı gibi birçok doğal dil işleme probleminde en iyi başarıyı göstermiştir. Buna rağmen, alana özgü çalışmalarda, gerek veri toplamanın zorluğu, zaman zaman da veri mahremiyeti sebebiyle ön-eğitimden (pre-training) geçirilmiş model bulmak zor olabilmektedir. Bu sebeple, Alsentzer vd. [46], klinik alana özgü eğitilmiş ClinicalBERT çalışmasını 2019 yılında sunmuştur. Çalışma kapsamında eğitilen modeller, BERT-Base [44] ve ince ayar yapılmış, BioBERT [47] modelleri üzerinden gerçekleştirilmiştir. Veri seti olarak, MIMIC-III içerisindeki klinik notlar seçilmiş, ve iki farklı deney yapısı kurgulanmıştır. Yapılan ilk deneylerde sadece hastaya ait taburcu notları kullanılırken (Discharge Summary BERT), ikinci deney kurgusunda bütün klinik notlar kullanılarak model (Clinical BERT) eğitimleri gerçekleştirilmiştir. Önerilen modeller, iki tane klinik VİT problemi (i2b2 2010 [108], i2b2 2012 [109]), bir tane medikal doğal dil çıkarım problemi (MedNLI [110]) ve iki farklı kimlik gizleme (de-identification) (i2b2 2006 [111] ve i2b2 2014 [112, 113]) problemi üzerinde test edilmiştir.

Çalışma kapsamında klinik notlar üzerinden eğitilen iki BERT modelinden bir tanesi, BERT-Base ile başlatılmış ve Clinical BERT ismi verilmişken, diğer model BioBERT ile başlatılmış, ve Clinical BioBERT ismi verilmiştir. İnce ayar yapılacak görevler (downstream tasks) kapsamında, BERT modellerinin ince-ayar yapılmasına izin verilmiştir. Deneyler esnasında yapılan veri ön işleme işlemi ve bu modellerin eğitim süresi, GeForce GTX TITAN X 12GB'lık ekran kartında yaklaşık olarak 17-18 gün sürmüştür. Sonuçlar incelendiğinde, önerilen modeller VİT ve MedNLI görevlerinde daha başarılı sonuçlar göstermiş, kimlik gizleme problemlerinde ise BioBERT çalışmasından kesin (exact) F1 skoruna göre %0.3 düşük performans göstermiştir. Yazarlar tarafından bu durumun sebebi, kimlik gizleme problemlerinde veri dağılım farklılığı olabileceği yönünde açıklanmıştır.

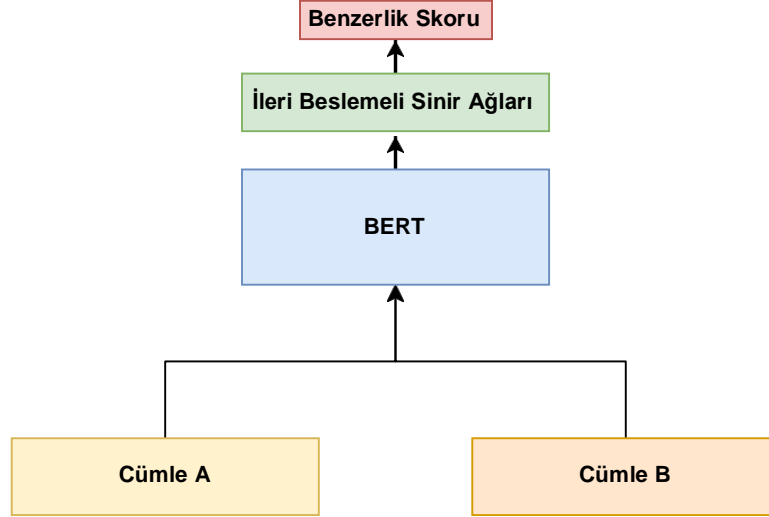
Tez kapsamında yapılan deneylerde klinik notların temsil edilebilmesi adına Alsentzer vd. yayınlamış olduğu ClinicalBERT modeli kullanılmıştır. İlgili model ve çalışma kodlarına buradan ³ erişilebilmektedir.

4.2.6 Sentence-BERT

BERT modeli kelime temsillerini öğrenmede oldukça iyi bir yöntem olmasına rağmen, içerisinde birden çok kelime bulunan bir cümleyi temsil etmek için bir takım ek yöntemler gerekmektedir. Bu problemi çözmek için ise literatürde temelde iki yöntem kullanılmaktadır. İlk ve genel olarak kullanılan yöntemlerden bir tanesi, kelime temsillerinin ortalamasını alarak cümle temsili elde etmektedir. İkinci olarak ise, özel [CLS] token vektör bilgisini, cümle temsili olarak kabul etmektedir. 2019 yılında Reimers ve Gurevych [95] tarafından yayınlanan Sentence-BERT çalışmasında, metin benzerliği görevi için cümle temsili elde etme noktasında, BERT yönteminin ürettiği kelime vektörlerinin ortalamasının alınmasının ve [CLS] vektörünün kullanılmasının geleneksel Glove yöntemine göre bile düşük performans verdiği gösterilmiştir. Buna ek olarak metin benzerliği görevi, orijinal BERT yöntemi ile çözülmeye çalışıldığında, bütün olası cümle ikililerinin Şekil 4.15'de görüldüğü üzere modele girdi olarak verilmesi gerekmektedir. Bu durum, n tane cümle olduğu düşünüldüğü durumda, $(n) \times (n-1) / 2$ kadar hesaplama maliyeti ($O(n^2)$) yaratmaktadır.

Bu durumu iyileştirmek adına önerilen bu çalışmada, içerisinde ikiz sinir ağları (siamese neural networks) kullanan Sentence-BERT (SBERT) isminde yöntem önerilmiştir. Bu yöntemde, üçlü hata fonksiyonu (triplet objective function) ve ikiz BERT modeli kullanılarak, bütün örnekler sadece bir kez modele girdi verilecek şekilde tasarlanmıştır. Eğitim aşamasında, Şekil 4.16'de görüleceği üzere, cümle A ve B özdeş iki BERT modeline girdi olarak verilmekte ve çıktıları ortalama-örnekleme yöntemi ile

³github.com/EmilyAlsentzer/clinicalBERT



Şekil 4.15: BERT modeli ile cümle benzerlik görevi mimarisi.

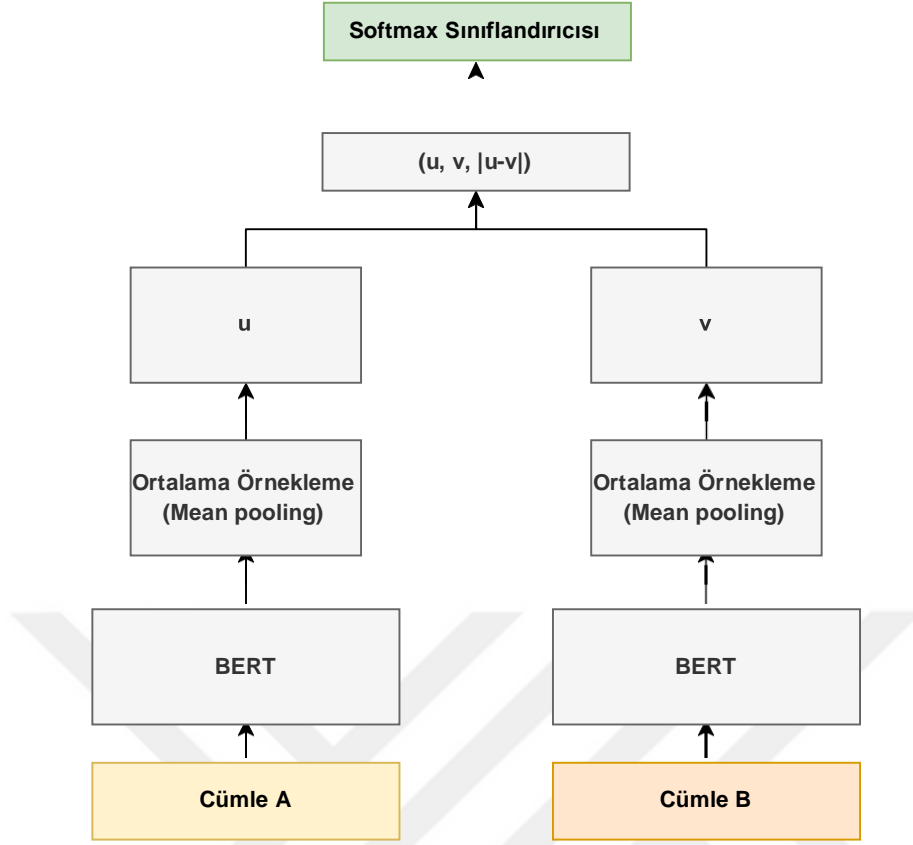
vektörlere dönüş, türülmektedir. Ardından A ve B cümle temsilleri birleştirilerek softmax sınıflandırıcısına gönderilmekte ve sonrasında ise eğitim gerçekleştirilmektedir. Bu sayede BERT modelinin aksine $O(n)$ 'lik bir hesaplama maliyeti ile görev tamamlanmış, olmaktadır.

Bölüm 6'de anlatılan çalışmada, klinik notlar gerekli ön işlemlerden geçirildikten sonra temsillerinin elde edilmesi için SBERT yöntemi kullanılmıştır. SBERT yöntemi sadece sonuçları almak için (inference-time) kullanılırken, mimari içerisindeki BERT modeli için de klinik alana özgü olarak eğitilmiş ClinicalBERT [46] kullanılmıştır.

4.3 Klinik Varlık İsimlerinin Tanımlanması

Doğal Dil İşleme alanının bir alt problem türü olan Varlık İsim Tanıma (VİT, Named Entity Recognition, NER), verilen bir metin içerisinde kişi, yer, organizasyon gibi bilgileri çıkartmayı amaçlamaktadır. Farklı dil ve problem türleri için eğitilmiş birçok VİT modeli literatürde bulunmaktadır [114]. Tez kapsamında ise medikal alandaki metinler içerisinde varlık isimlerinin çıkartılması amaçlandı için, literatürdeki sağlık alanında geliştirilen VİT modelleri incelenmiş, ve med7 [53] çalışması kullanılmıştır.

Kormilitzin vd. [53] tarafından önerilen med7 çalışması, girdi olarak verilen medikal metin içerisinde 7 farklı kategoriye sınıflandırabilmektedir. Bu kategoriler: ilaç ismi (drug names), ilaç alım bilgisi (route), ilaç kullanım sıklığı (frequency), doz bilgisi (dosage), ilaç gücü (strength), ilaç formu (form) ve ilaca devam etme süresidir (duration). Model ilk olarak öz-denetimli öğrenme (self-supervised learning) yöntemi ile bir sonraki kelimeyi tahmin etme problemi üzerinde ön-egitimden geçirilmiştir. Bu aşamada, bu tez kapsamında da kullanılan MIMIC-III veri seti içerisindeki klinik notlar kullanılmıştır.



Şekil 4.16: Sentence-BERT mimarisi.

MIMIC-III veri seti içerisinde yaklaşık olarak 2 milyon hasta notu bulunmaktadır. Bu ön-eğitim aşamasından sonra ise, modele VİT problemi üzerinde ince ayar yapılmıştır. VİT problemi için ince ayar yapılırken ise "2018 National NLP Clinical Challenge (n2c2)" yarışmasında yayınlanmış, etiketli VİT veri seti kullanılmıştır. Ayrıca veri/etiket sayısını arttırmak adına arastırıcılar, kural tabanlı, sözlük tabanlı, insanın da içinde bulunduğu çeşitli yöntemler ile veri sayısını arttırmışlardır. Model eğitimi sırasında evrimsel sinir ağları kullanılmış ve test veri seti üzerinde katı (strict) ortalama mikro F1 skorunda %89.3, rahat (lenient) ortalama mikro F1 skorunda ise %95.7 sonucu elde edilmiştir. Tez kapsamında, klinik notlar içerisinde medikal terimleri çıkartma problemi için ön-eğitimden geçirilmiş, kodu ve modeli ⁴ paylaşılmış med7 yöntemi kullanılmıştır. MIMIC-III veri seti içerisindeki bir klinik notun içerisindeki cümle üzerinde, med7 yönteminin vermiş olduğu sonuçlar Şekil 4.17'de gösterilmiştir.

4.4 İlaçların Temsili

Geleneksel yapay öğrenme problemlerindeki önemli bir zorluk, girdi olarak verilecek veri kümesini hazırlamaktır. Veri kümesinin hazırlanmasında, veri içerisindeki bilgi

⁴<https://github.com/kornilitsin/med7>

Agitation: For bouts of agitation, please consider a trial of **haldol DRUG** **1 mg STRENGTH** to **2.5 mg STRENGTH** (**PO ROUTE** or **IV ROUTE**)
or **zyprexa DRUG** **disintegrating DOSAGE** tablet **2.5 mg STRENGTH** to **5 mg STRENGTH** (placed under the tongue)

Şekil 4.17: med7 çalışması kapsamında e ğitilen klinik VİT modelinin, MIMIC-III

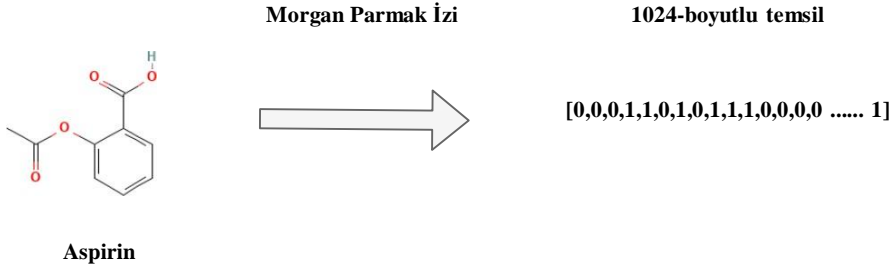
içerisindeki bir klinik not cümlesi üzerindeki örnek çıktısı.

kaybedilmemeye çalışılırken, bir yandan da veri boyutunun, modellere verilebilecek düzeye indirgeme işlemi bir arada gerçekleştirilir. Geleneksel yapay öğrenme problemlerinin aksine, derin öğrenme yönteminde ham veri de doğrudan modellere girdi olarak verilebilmektedir. Buna rağmen, moleküler veri ile yapılan çalışmalarda, moleküler verinin hem yapay öğrenme hem de derin öğrenme algoritmalarına girdi olarak en başarılı şekilde nasıl verilebilece ği hala tartışılabilir bir noktadır.

Bölüm 7’de yapılan çalışmada hastalara ait ilaç bilgileri, ilaçların farklı moleküler temsilleri kullanılarak önerilen çok-kipli derin öğrenme tabanlı modele verilmiştir. Bu sebeple, veri seti içerisinde hastalara ait klinik ilaç isimlerinin moleküler temsillere dönüşürülmesi gerekmiştir. Bu sebeple, farklı çalışma prensiplerine sahip ECFP, MACCS, Mol2Vec ve Smiles-Transformer yöntemleri uygulanmıştır. Yöntemlerin detayları bu bölüm altında açıklanmıştır.

4.4.1 Extended-Connectivity Fingerprints (ECFP)

Genişletilmiş Bağlantı Parmak İzi (Extended-Connectivity Fingerprints, ECFP), 2010 yılında Rogers vd. [58] tarafından yayınlanan ve kural tabanlı olarak hesaplanan bir moleküler parmak izi yöntemidir. Çok ölçekli moleküler alt yapıları, tam sayılara sifreleyerek sabit uzunlukta ikili (sıfır birden oluşan) vektörler oluşturmaktadır. Vektörler içerisinde yer alan 1 değeri moleküldeki belirli alt yapıların varlığını, 0 ise yokluğunu temsil etmektedir. Bu dizi içerisindeki i . tamsayı, ilişkili olduğu atom hakkındaki bilgilerin yanı sıra, o atomun i bağındaki atomlar ve bağlar hakkındaki bilgileri, yani o atomun i bağlarındaki bileşimin alt yapısını kodlamaktadır. Bu yapı, ECFP yönteminin hesaplamalı kimyada kullanımının birçok problem için yararlı olmaktadır. Bu yapının bir diğer avantajı, daha önceden görmediği herhangi bir molekül içinde yeni yapısal sınıfları temsil edebilmesidir. Aynı zamanda ECFP yönteminin hızlı hesaplanabilir ve kolay yorumlanabilir olması, literatürdeki birçok çalışmanın ECFP yöntemine deneylerinde yer vermesine sebep olmaktadır. Morgan algoritmasının [115] bir varyasyonunu kullanan bu yöntem, bir bileşikteki tüm tanımlanmış alt yapılar hakkında bilgi toplayarak, girdi bileşiğinin boyutundan bağımsız olarak sabit uzunlukta bir vektör üretmektedir (1024 veya 2048 boyutlu). Şekil 4.18’de ECFP için örnek bir gösterim yapılmıştır.



Şekil 4.18: ECFP yöntemi örnek gösterimi.

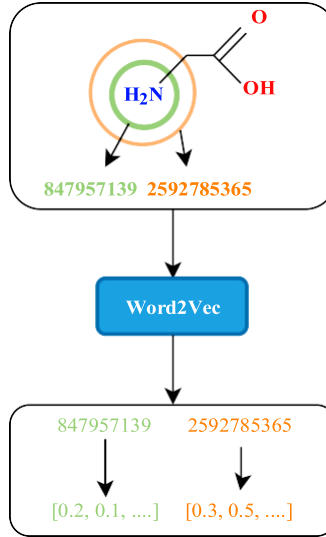
4.4.2 Molecular Access System (MACCS)

Tez kapsamında kullanılan ilaçların molekül karakterizasyonu için kullanılan bir diğer yöntem ise yapısal anahtar tabanlı parmak izlerine ait olan MACCS yöntemidir. Zaman zaman MDL (Molecular Design Limited) anahtarı olarak da isimlendirilen MACCS [59], literatürde yaygın olarak kullanılan parmak izlerinden bir tanesi olup, 2002 yılında Durant vd. tarafından önerilmiştir. Basit bir tasarıma sahip olan MACCS parmak izleri, çeşitli atomları içeren moleküler alt yapılarla beraber, 166 anahtara sahip sabit uzunlukta vektörel temsiller öğrenmektedir. MACCS vektörü içerisindeki bitler, molekül içerisindeki ilgili yapısal parçaya atanırken, bu değerlerin 1 olması bileşiğin varlığını, 0 olması ise bileşiğin yokluğunu ifade etmektedir. Yapısal anahtar yöntemine dayanan MACCS parmak izinin ECFP yönteminden farkı ise, özelliklerin önceden tanımlanmış, belirli altyapısını temsil etmek için bir dizi kullanmasıdır. Bu çalışmada kapsamında, MACCS gösterimlerini hesaplayabilmek adına ise DeepChem [82] kullanılmıştır.

4.4.3 Mol2Vec

Molekül temsillerini öğrenmek için önerilen Mol2Vec çalışması [60] Jaeger vd. tarafından 2018 yılında önerilmiştir. Doğal dil işleme yöntemlerinden esinlenerek geliştirilen bu yöntemde, moleküllerin vektörel temsilleri gözetimsiz yapay öğrenme yöntemi ile öğrenilmiştir. Birbirine yakın kelimelerin vektörel temsillerini birbirine yaklaştıran ve Bölüm 4.2.1’de detayları anlatılan Word2Vec yöntemine benzer şekilde, Mol2Vec yöntemi de birbirine yakın olan moleküler altyapı bilgisini kullanarak moleküler temsilleri öğrenmektedir.

Modelin ön-eğitimi (pre-training) gerçekleştirmek için ZINCv15 [116] ve ChEMBL v23 [117] veri setleri kullanılmıştır. Bu iki veri tabanından gelen veriler birleştirilmiş, ve tekrarlı veriler silinmiştir. Daha sonra ise çeşitli diğer kriterlere göre elemeler yapılmıştır. Son aşamada ise, Mol2Vec eğitiminde kullanılmak üzere 19.9 milyon bileşik bilgisi elde edilmiştir. 19.9 milyon bileşik bilgisi iki katmanlı bir yapay sinir ağı olan Word2Vec yöntemi ile eğitilmiştir. Şekil 4.19’de gösterimi olan bu yöntem, girdi olarak aldığı bir

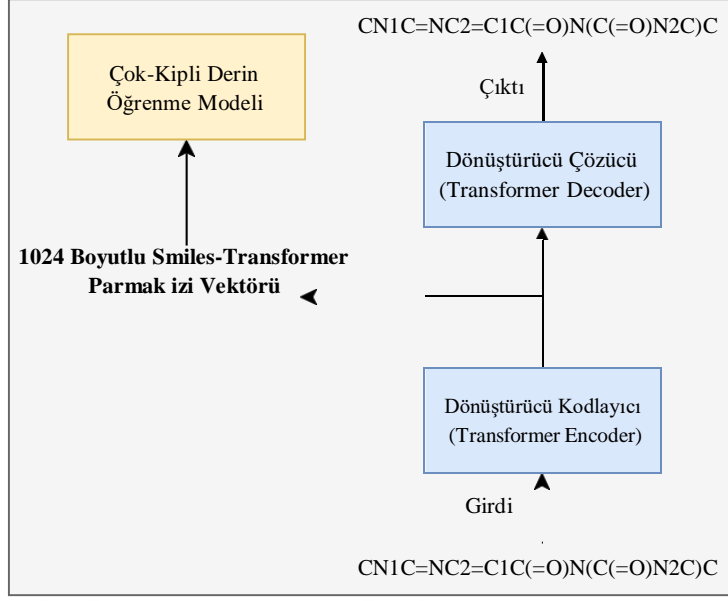


Şekil 4.19: Mol2Vec yöntemi örnek gösterimi.

bileşiğin 100 ve 300 boyuttaki temsillerini yaratabilmeyi öğrenmiştir.

4.4.4 Smiles-Transformer

Geleneksel olarak, moleküler parmak izleri kural tabanlı algoritmalar ile hesaplanmasına rağmen, bu hesaplama yönteminin küçük boyutlu veri kümeleri için düşük başarımlar gösterdiği Honda vd. [61] tarafından 2019 yılında önerilen çalışmada belirtilmiştir. Bu çalışmada, moleküler parmak izlerini geleneksel yöntemler ile oluşturmak yerine dönüştürücü (transformer) tabanlı algoritmadan esinlenerek oluşturulmuştur. Basitleştirilmiş Moleküler Giriş Sistemi (Simplified Molecular Input Line Entry System, SMILES) moleküler yapıların tabanlı olarak tanımlayan bir yöntemdir. Bu çalışmada, dönüştürücü mimarilerde olduğu gibi etiketsiz veri ile gözetimsiz olarak ön eğitim işlemi gerçekleştirilmiştir. Önerilen kodlayıcı-çözücü mimarisine 4 dönüştürücü bloğu yerleştirilmiştir. Her dönüştürücü blok içerisinde, 4-başlı dikkat mekanizması, 256 boyutlu temsil vektör girişi ve 2 tane lineer katman bulunmaktadır. SMILES-Transformer yönteminin ön eğitimi (pre-training), ChEMBL24 içerisinde rastgele seçilen 861,000 etiketsiz SMILES ile gerçekleştirilmiştir. SMILES'lar modele girdi olarak verilebilmesi için sembollere ayrıştırılmış ve sıfır-kodlama yöntemi ile vektörize edilmiştir. Önerilen ağ modeli, 5 iterasyon boyunca eğitilmiş ve çapraz entropi hata fonksiyonu minimize edilmeye çalışılmıştır. Eğitim aşamasından sonra parmak izlerini çıkartmak için girdi olarak verilen her molekülün 1024-boyutlu parmak izi alınabilir olmaktadır. Detayları anlatılan bu modelin özet gösterimi Şekil 4.20'de gerçekleştirilmiştir.



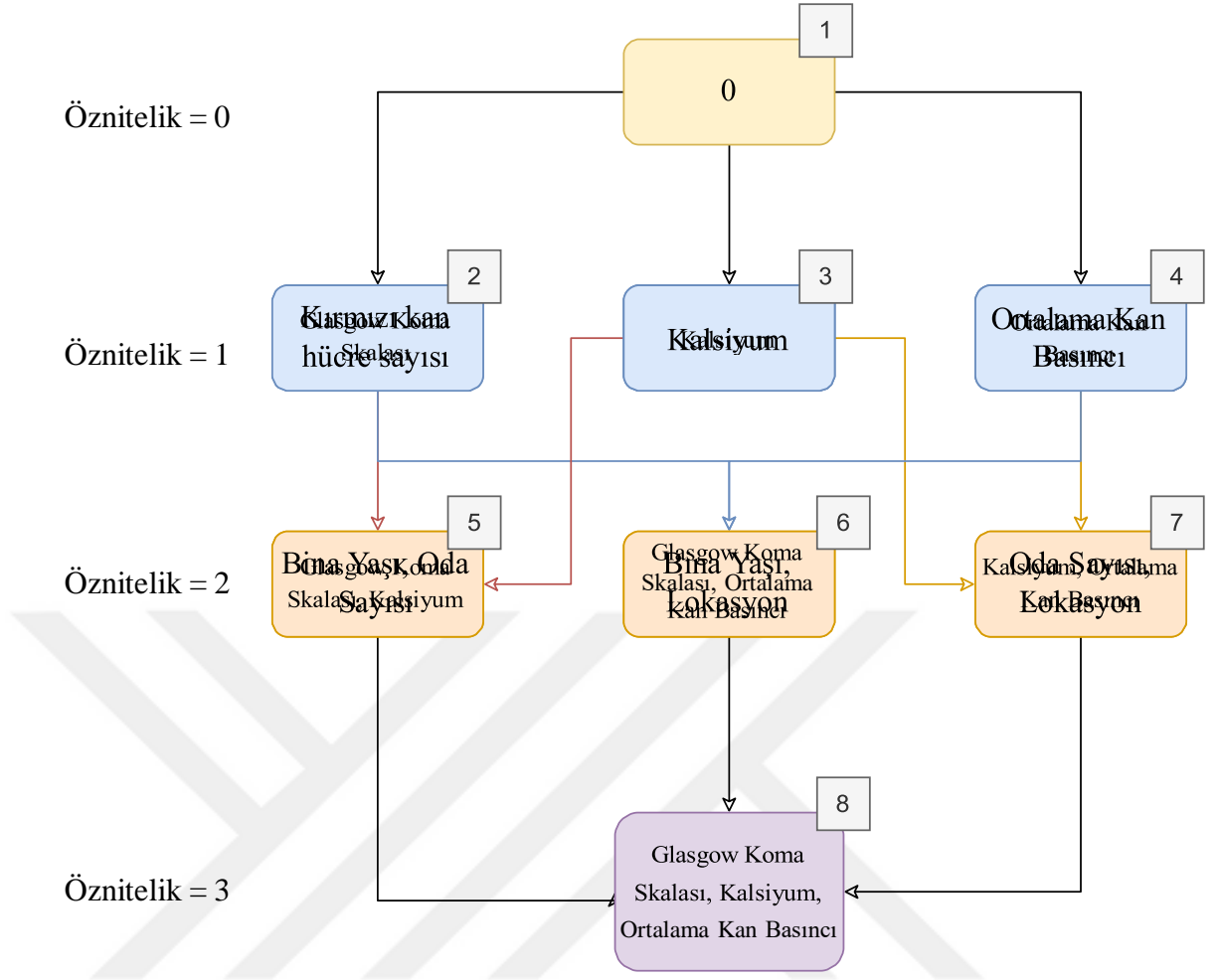
Şekil 4.20: SMILES-Transformer yöntemi ön-eğitim modeli.

4.5 SHapley Additive exPlanations (SHAP)

SHapley Additive exPlanations (SHAP) yöntemi 2017 yılında Lundber ve Lee [75] tarafından modellerin açıklanabilirliğini arttırmak amacıyla önerilmiştir. Bu yöntem sayesinde modelden bağımsız bir şekilde, modelin tahmin karar verirken hangi özneliklere dayanarak bu tahmini yapabildiği tespit edilebilir olmuştur. SHAP yöntemi, temelinde oyun teorisindeki Shapley değerlerini kullanmaktadır. Oyun teorisinde temelde en az iki nesnenin bulunması gerekir: oyuncular ve oyun. Oyun teorisi içerisindeki Shapley değerleri aslında her bir oyuncunun oyuna olan katkısını hesaplamaktadır. Bu fikri, yapay öğrenme problemlerine uyguladığımızda ise oyun, herhangi bir örnek verildiğinde o örneğe ait model çıktısına, oyuncular ise modelin özneliklerine denk gelmektedir. Bu sayede SHAP, her bir özneliğin model çıktısına katkısını inceleyerek modelleri açıklanabilir hale getirmektedir.

Shapley değerleri, tek bir özneliğin önemini belirlemek için özneliklerin her olası kombinasyonunun sonucunun dikkate alınması gerektiği fikrine dayanır. Bu sebeple SHAP, eğer öznelik sayısı F olarak kabul edilirse, bütün öznelik kombinasyonlarıyla, yani 2^F sayı kadar model e ğitir. Örneğin; bir hastanın yoğun bakımda kalma süresini tahmin etmek için geliştirilecek olan modelde, "Glaskow Koma Skalası", "Kalsiyum" ve "Ortalama Kan Basıncı" isminde üç öznelik olduğunu varsayalım. O halde, Şekil 4.21'de görüldüğü üzere $2^3 = 8$ adet model e ğitimi gerekmektedir.

Şekil 4.21'de gösterildiği üzere, dikdörtgen (node) içerisindeki ifadeler modellerin hangi öznelik ile eğitildiğini temsil ederken, kutular birbirine bağlayan çizgiler (edge) ise fazladan eklenen özneliğin marjinal katkısını ifade etmektedir. Şekil içerisinde



Şekil 4.21: SHAP öz nitelik kombinasyon gösterimi.

kutuların sağ üst kısmında verilen numaralar, ilgili deneyin kaçınıcı deney oldu ğunu ifade etmektedir. 1 numaralı deney, herhangi bir öz nitelik kullanılmadan gelen girdinin ortalamasının dönüldüğü deneydir. 2 numaralı deneyde ise, model içerisinde sadece "Glasgow Koma Skalası(GKS)" öz niteli ği kullanılarak eğitim yapılmıştır. Bu durumlar ışığında, 2.deneydeki "Glasgow Koma Skalası"nın modele olan marjinal katkısı(MK) x_0 örne ği için 1. ve 2. modellerinin tahmin farkı olarak ifade edilebilmektedir.

$$MK_{2, \text{GKS}, x_0} = 2. \text{ Deney Model Tahmini} - 1. \text{ Deney Model Tahmini} \quad (4.14)$$

Glasgow Koma Skalasının toplam marjinal faydasını hesaplayabilmek için ise 4 farklı modelin çıktılarının karşılaştırılmasından hesaplanacak marjinal faydanın hesaplanması gerekmektedir. Bu deneylerin; MK_2 için (1. ve 2. deney), MK_5 için (3. ve 5. deney), MK_6 için (4. ve 6. deney), ve MK_8 için (7. ve 8. deney) arasında yapılması gerekmektedir. Bu deneylerden hesaplanacak olan marjinal katkıların, gerekli ağırlıklar ile çarpılması ile, Glasgow Koma Skalasının, x_0 örne ği için SHAP değerini hesaplamaktadır. Bu

işlemlerin detaylarına şa ğıdaki formülde paylaşılmıştır.

$$\text{SHAP}_{\text{glasgow koma skalası, } x_0} = w_1 \diamond MK_2 + w_2 \diamond MK_5 + w_3 \diamond MK_6 + w_4 \diamond MK_8 \quad (4.15)$$

Bu işlemden sonra ise yukarıdaki formülde de belirtilen w_1, w_2, w_3, w_4 ağırlıkları hesaplanmalıdır ($w_1 + w_2 + w_3 + w_4 = 1$ olmalıdır.). Bu hesap esnasında SHAP iki tane varsayım ile ilerlemektedir. İlk olarak, bir öznitelikle eğitilen modellere yapılan tüm marjinal katkıların ağırlıklarının toplamının, iki öznitelikle eğitilen modellere yapılan tüm marjinal katkıların ağırlıklarına eşit olmasıdır. Bu diğ er sayıdaki öznitelikle eğitilen modeller için de geçerlidir. Bu varsayımdan, üzerinde ilerlediğimiz örnek için $w_1 = w_2 + w_3 = w_4$ şeklinde bir eşitlik doğ maktadır. İkinci olarak ise, aynı öznitelik sayısı ile eğitilen bütün modellerin ağırlıklarının birbirine eşit olma koşuludur. Bu sebeple, $w_2 = w_3$ şeklinde bir başka eşitlik oluş maktadır. Bu kısıtlar beraber düşün üldüğ ünde, f öznitelik ile eğitilen modellerin tüm marjinal katkılarının sayısı:

$$f \diamond \begin{matrix} \rightarrow \\ F \\ \rightarrow \\ f \end{matrix} \quad (4.16)$$

olarak ifade edilebilir. Yukardaki örneğ e ait w_1, w_2, w_3, w_4 değ erlerini hesaplamak için, modele ait öznitelik sayısı f ve o öznitelik sayısı ile eğ itilen model sayısı, F ile ifade edildiğ inde, ağırlık sonuçlarına şa ğıdaki gibi çıkmaktadır.

$$w_1 = [1 \diamond \begin{matrix} \rightarrow \\ 3 \\ \rightarrow \\ 1 \end{matrix}]^{-1} = 1/3 \quad (4.17)$$

$$w_2 = [2 \diamond \begin{matrix} \rightarrow \\ 3 \\ \rightarrow \\ 2 \end{matrix}]^{-1} = 1/6 \quad (4.18)$$

$$w_3 = [2 \diamond \begin{matrix} \rightarrow \\ 3 \\ \rightarrow \\ 2 \end{matrix}]^{-1} = 1/6 \quad (4.19)$$

$$w_4 = [3 \diamond \begin{matrix} \rightarrow \\ 3 \\ \rightarrow \\ 3 \end{matrix}]^{-1} = 1/3 \quad (4.20)$$

Sonuçlar incelendiğ inde, $w_1 = w_2 + w_3 = w_4$ eşitliğ inin, $w_2 = w_3$ eşitliğ inin ve $w_1 + w_2 + w_3 + w_4 = 1$ eşitli ğ inin korunduğ u gözlemlenmektedir.

SHAP yönteminin kolay anlaş ılabilir olması için üç örnek üzerinden verilen bu örnek genelleş tirme ve daha resmi olarak ifade etmek istenirse ş u şekilde açıklanabilir. Orijinal SHAP çalış masında da [75] açıklandığı gibi, F özniteliğ e sahip bir modelin tüm öznitelik alt kümeleri olan S üzerinde eğ itilmesi gerekmektedir ($S \hat{=} F$). Böylece, her bir öznitelik, o özniteliğ in dahil edilmesinin model tahmini üzerindeki etkisini

temsil edecek bir önem değeri atanmış olur. Bu değeri hesaplamak için ise, örneğin i . öznelik için, $f_{S \setminus \{i\}}$ ve f_S modelleri eğitilir. Daha sonra ise her iki modelin tahminlerinin farkları, $f_{S \setminus \{i\}} - f_S(S)$, hesaplanır. Bir özneliği çıkartmanın etkisi modeldeki diğer özneliklere de bağlı oldu ğundan, önceki farklar tüm olası $S \in \mathcal{F} \setminus \{i\}$ alt kümeler için hesaplanmaktadır. O halde, herhangi bir öznelik i için SHAP değeri aşağıdaki formüle göre hesaplanabilmektedir.

$$\phi_i = \sum_{S \in \mathcal{F} \setminus \{i\}} \frac{|S|!(|F|-|S|-1)!}{|F|!} (f_{S \setminus \{i\}} - f_S(S)) \cdot \bar{\phi} \quad (4.21)$$

Bu yöntem ile, bütün örnek veriler için, model içerisindeki bütün özneliklerin SHAP değerleri bulunabilmektedir. Bölüm 7’de SHAP yöntemi uygulanarak, özneliklerin önemleri bulunmaya çalışılmıştır.

4.6 Performans Metrikleri

Tez kapsamında üzerinde çalışılan dört klinik problem de (hastane içinde mortalite, yoğun bakım içinde mortalite, YBÜ >3, YBÜ >7) ikili sınıflandırma problemi olarak tasarlanmıştır. Bu sebeple, eğitilen modeller sınıflandırma problemlerinde kullanılan metriklerden yararlanılarak karşılaştırılmış, ve sonuçlar okuyucu ile paylaşılmıştır.

Şekil 4.22’de görülen karışıklık matrisi üzerinden farklı sınıflandırma metrikleri hesaplanabilmektedir. Bu tez kapsamında geliştirilen model sonuçları üç farklı metrik ile değerlendirilmiştir. Bu metrikler, F1 skoru, Receiver Operating Characteristics (ROC) altında kalan alan metriği (Area under ROC curve, AUROC) ve Kesinlik (Precision)-Duyarlılık (Recall) eğrisi altında kalan alan (Area under Precision-Recall Curve, AUPRC) metriği olarak seçilmiştir.

F1 skoru, Kesinlik 4.22 (Precision) ve Duyarlılık 4.23 (Recall) metriklerinin harmonik

	Gerçek Değerler	
Tahmini Değerler	Doğru Pozitif (TP)	Yanlış Pozitif (FP)
	Yanlış Negatif (FN)	Doğru Negatif (TN)

Şekil 4.22: Karışıklık matrisi.

ortalaması ile hesaplanmaktadır. F1 metriği içerisinde kullanılan Kesinlik ve Duyarlılık metriklerinin formülleri aşağıda paylaşılmıştır.

$$\text{Kesinlik} = \frac{TP}{TP + FP} \quad (4.22)$$

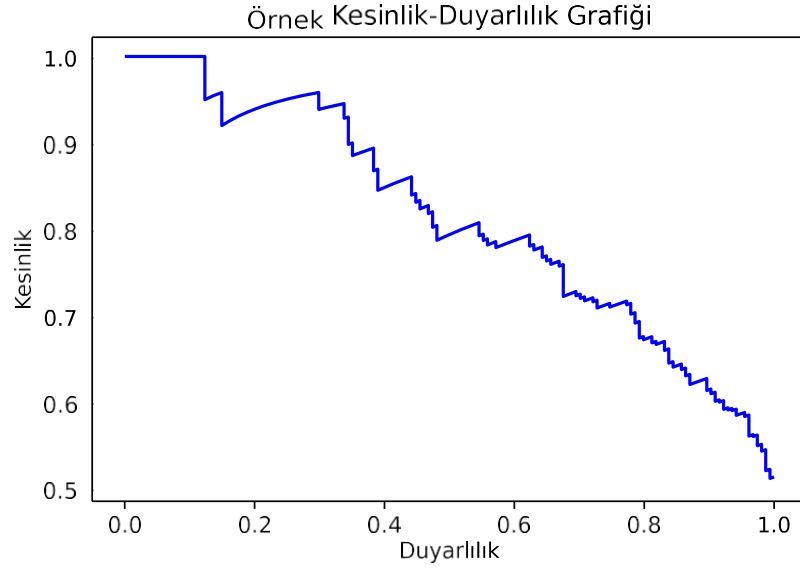
$$\text{Duyarlılık} = \frac{TP}{TP + FN} \quad (4.23)$$

Kesinlik metriği, modelin pozitif olarak tahminlediği örneklerden kaç tanesinin pozitif olduğunu hesaplayarak modelin kesinliği ile ilgili bilgi verirken, Duyarlılık metriği, modelin gerçekte pozitif olan örneklerden kaçınıcı şekilde pozitif olarak

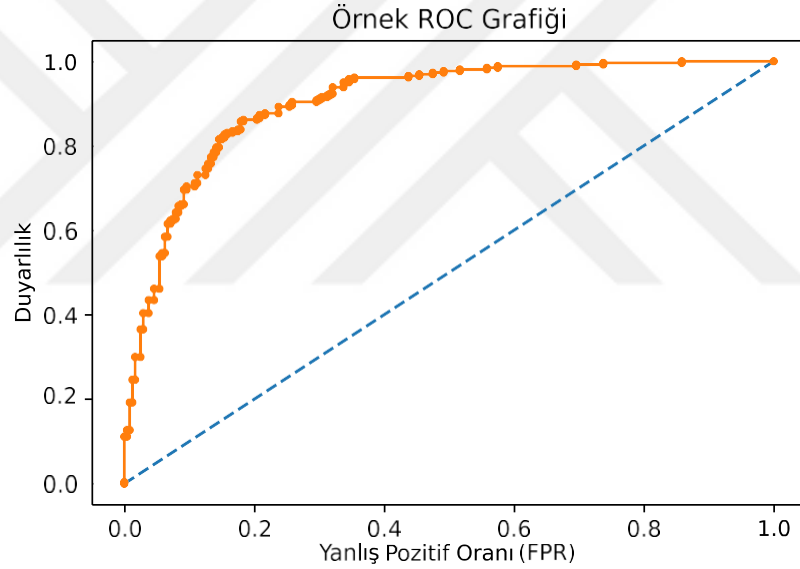
tahmin edebildiğini göstermektedir. Hedeflenen amaca göre bu iki metrikten birisi kullanılabilir gibi, Kesinlik ve Duyarlılık metriklerinin harmonik ortalaması alınarak F1 skoru 4.24 kullanılabilir.

$$F_1 = \frac{2 \times \text{Kesinlik} \times \text{Duyarlılık}}{\text{Kesinlik} + \text{Duyarlılık}} \quad (4.24)$$

AUPRC (Area Under Precision Recall Curve) özellikle sınıf dağılımı sorunu (class imbalance) olan problemler için kullanılan bir performans metriğidir. Mortalite ve YBÜ’de kalma süresini tahmin etme problemleri de birçok klinik problem de olduğu gibi sınıf dağılımı sorunu içerdiğinden model başarımları bu metrik ile de ölçülerek karşılaştırılmıştır. AUPRC metriği, Kesinlik-Duyarlılık (PR) grafiği altında kalan alan hesaplanarak bulunmaktadır. PR grafiği, Kesinlik ve Duyarlılık arasındaki ödünleşmeyi farklı değerler için ölçerek oluşturulmaktadır. Şekil 4.23’de örnek bir PR grafiği gösterilmiştir. PR grafiklerinin hesaplanması sol üst köşeden (eşik değeri = 1 seçilerek) başlar ve sağ alt köşede (eşik değeri = 0) bitmektedir. Grafiğin altında kalan alan ise AUPRC değerini vermektedir. Bu alanın hesaplanması için farklı yöntemler olmakla beraber genellikle ortalama Kesinlik skoru kullanılarak bu değer hesaplanabilmektedir. Model başarımlarını karşılaştırmak için kullanılan son metrik ise AUROC (Area Under Receiver Operating Characteristics Curve) metriğidir. ROC, Şekil 4.24’de görüldüğü gibi x-koordinatında Yanlış Pozitif Oranı, y-koordinatında da Duyarlılığı bulunduran bir grafikdir. AUROC değeri ayrımcı (discriminative) bir performans metriği olmakla beraber modelin pozitif ve negatif örnekler arasında ayırım yapma kabiliyeti hakkında bilgi vermektedir. ROC grafiği, hesaplanmaya her zaman sol alt köşeden başlamakta (eşik değeri = 1 seçilerek) ve her zaman sağ üst köşede sonlanmaktadır (eşik değeri = 0). ROC grafiğinin altında kalan alan ise hesaplanarak AUROC değeri elde edilmektedir.



Şekil 4.23: Örnek Kararlılık-Duyarlılık (PR) grafiği.



Şekil 4.24: Örnek ROC grafiği.

4.7 Kullanılan Kütüphaneler ve Çalışma Ortamı

Tez kapsamında gerçekleştirilen deney ve çalışmalarda kullanılan kütüphaneler, kelime temsil modelleri ve ortam bilgileri bu bölümde okuyucu ile paylaşılmıştır. Yapılan deneylerde kullanılan veri seti olan MIMIC-III'e erişim Bölüm 3'de detaylandırılmıştır. MIMIC-III veri setini saklamak için gereken alan minimum 40GB civarındadır. Veri setine erişim aşamasından sonra, zaman serisi öznitelikleri yaratabilmek ve veriyi ön işlemeden geçirebilmek için detaylar Bölüm 3.2'de anlatılan MIMIC-Extract çalışması

kullanılmıştır. MIMIC-Extract çalışmasının tekrarlanabilmesi için, çalışma yapılan ortamda PostgreSQL 9.4 ve üzeri veritabanı ve conda paket yöneticisinin yüklü olması gerekmektedir. MIMIC-Extract çalışmasının tekrarlanacağı bilgisayarda minimum 50GB Rastgele Erişimli Hafıza (RAM) bulunması beklenmektedir. Yapılan diğer veri ön işleme aşamalarında ise veriyi okumak ve veri üzerinde çeşitli işlemler gerçekleştirmek için Pandas [118], matris işlemleri için Numpy [119] kütüphanelerinden yararlanılmıştır. Verinin eğitim, validasyon, test kümesine ayrılması, normalizasyon işlemleri ve modellerin başarımlarının hesaplanması gibi işlemler için ise Scikit-learn [120] kütüphanesi kullanılmıştır.

Bölüm 5 ve Bölüm 6’de yapılan çalışmalarda hastalara ait klinik notlar modellere girdi olarak verilmiştir. Klinik notların vektörel temsillerinin elde edilebilmesi için önceden eğitilmiş Word2Vec [48] ve FastText [48] kullanılmıştır. Bu yöntemler Gensim [121] kütüphanesi aracılığıyla modellerde kullanılmıştır. Klinik notların doğrudan temsili için kullanılan ClinicalBERT modeline bu bağlantıdan ulaşılmış ve modelin kullanılabilmesi için Sentence-BERT [95] kütüphanesinden faydalanılmıştır. Önerilen derin öğrenme tabanlı modeller ise Tensorflow [122] ve Keras [123] kütüphaneleri ile eğitilmiştir. Klinik notlar içerisinde medikal terimlerin çıkartılabilmesi için ise, medikal varlık isim tanıma modeli olan ve arka planda spacy [124] kütüphanesinden yararlanan med7 [53] çalışması kullanılmıştır.

Bölüm 7 kapsamında yapılan çalışma ve ilgili deneylerde, hastaya ait ilaç bilgilerinin moleküler yapıları vektörel temsillere dönüştürülmüştür. Bu aşamada, ilaç isimlerinin SMILES bilgisine dönüştürülmesi için pubchempy⁵ kütüphanesi sayesinde Pubchem sitesinden bilgilere erişilebilmiştir. Ardından SMILES formatındaki verilerin ECFP, MACCS, Mol2Vec formatlarına dönüştürülebilmeleri için Rdkit⁶ ve DeepChem [82] kütüphanelerinden yararlanılmıştır.

⁵<https://pubchempy.readthedocs.io/en/latest/>

⁶<https://www.rdkit.org>

5. ZAMAN SERİSİ ve MEDİKAL TERİMLER ile TAHMİN ETME

5.1 Motivasyon

Yoğun bakımda yatan bir hastanın klinik ölçüm değerlerini, laboratuvar test sonuçlarını ve diğer verilerini kullanarak, hastanın genel sağlık durumunu anlayabilmek önemli bir klinik problemdir. Bu çalışmada, bir önceki çalışmada [70] olduğu gibi mortalite (hastanede ve YBÜ'de) ve yoğun bakımda kalma süresi (3 ve 7 günden büyük) klinik problemleri üzerinde çalışılmıştır. Her iki problemin çıktısında hastalara uygulanacak tedavi yöntemlerini belirleme, insan hayatını kurtarma ve hastane kaynaklarının planlanması için oldukça önemlidir. Literatürde klinik olayları tahmin etmeye yönelik yapılan çalışmalarda genellikle hastaya ait hastalık kodları (ICD kodları), laboratuvar sonuçları, gözlem verileri gibi yapısal veriler kullanılmış, ve yapısal olmayan ESK verilerinden çok fazla yararlanılmamıştır. ESK verileri hasta hakkında geçmiş, yapısal verileri içermekle beraber doktor, hemşire, radyoloji uzmanı ve diğer klinik uzmanların hastalar için yazdığı klinik notlarında içermektedir. Metin verilerinin işlenmesinin yapısal verilere göre daha zor olması ve klinik metinlerin kendine ait jargonunun bulunması sebebiyle, klinik notların kullanımı, yapısal verilerin kullanımına nazaran daha az tercih edilmiştir. Bu zorluklara rağmen, klinik notların içerisindeki zengin içeriği kullanmayı hedefleyen çalışmalar ortaya çıkmıştır [70, 65]. Hastaya ait laboratuvar sonuçları, gözlem verileri gibi yapısal verilerin yanında klinik notlar da yapay öğrenme modellerine girdi olarak verilerek üzerinde çalışılan klinik problemlerin başarımları artırılmaya çalışılmaktadır. Klinik notların genellikle uzun, dil bilgisi olarak yanlış, veya içerisinde gereksiz bilgiler bulundurabilmesinden dolayı klinik notların ham halini modellerde kullanmak zorlu bir süreç gerektirebilmektedir. Bölüm 4.3'de detayları anlatılan Varlık İsim Tanıma (NER) yönteminin klinik alana özelleştirilmiş modeli ile klinik notlar içerisinde bulunan medikal varlık isimleri çıkartılarak tahmin modellerine girdi olarak verilmesinin klinik problemlerde başarımları artırabileceği öngörülmektedir.

Bu çalışmanın literatüre katkılarına aşağıdaki gibi sıralanabilir:

- Klinik notlar içerisinden çıkartılan medikal terimlerin nasıl temsil edileceği önemli bir konudur. Farklı kelime temsil yöntemleri, yapıları itibarıyla aynı kelimenin farklı semantik ve sözdizimsel özelliklerini yakalayabilmektedir. Bu

çalışmada, Word2Vec, FastText, ve bu temsillerin birleştirilmiş hali (Word2Vec + FastText) yöntemleri ile deneyler yapılmış, ayrıca bu yöntemlerin başarıları birbirleri ile karşılaştırılmıştır.

- Medikal terimlerin efektif bir şekilde temsil edilebilmesi için vektörel ortalama alma, Doc2Vec, 1D CNN gibi farklı yöntemler denenmiştir. Yapılan deney sonuçları incelendiğinde 1D CNN yönteminin en başarılı sonucu verdiği görülmüştür.
- Hastanede mortalite, YBÜ'de mortalite, YBÜ'de 3 günden fazla kalma ve YBÜ'de 7 günden fazla kalma problemleri için özgün, tekrarlanabilir, derin öğrenme tabanlı çok-kipli model önerilmiştir.

5.2 Önerilen Yöntem

Çalışma kapsamında eğitilen bütün modeller MIMIC-III veri seti üzerinde gerçekleştirilmiştir. MIMIC-III veri seti içerisinde yapay öğrenme modellerinde girdi olarak kullanılacak zamana bağlı özelliklerin çıkarımı için Bölüm 3.2'de detaylandırılan MIMIC-Extract çalışması uygulanmıştır. Bütün deneylerde, hastanın yoğun bakıma yatışından itibaren ilk 24 saatlik zaman dilimi içerisinde toplanan verileri kullanılmıştır. Aşağıdaki kriterler uygulanarak çeşitli klinik notlar ve hasta verileri derlem içerisinde çıkartılmıştır. Bu kriterler:

- Veri sızıntısını engellemek adına kategorisi "hasta taburcusu" olan klinik notlar derlem içerisinde çıkartılmıştır. (Hasta taburcusu notları içerisinde hastanın hastaneye girdiği, çıktığı zaman bilgileri veya mortalite olup olmadığı bilgisi bulunabilmektedir.)
- Zaman damgası bulunmayan klinik notlar derlem içerisinde çıkartılmıştır.
- Kalan klinik notlar içerisinde YBÜ'de yatan hastaların sadece ilk 24 saat içerisinde değerlendirilen klinik notları derlem içerisinde tutulmuştur.
- Bu adımlardan sonra hiç klinik notu bulunmayan veya klinik notları içerisinde hiç medikal terim çıkartılmamış, hastalar derlem içerisinde çıkartılmıştır.

Belirtilen kriterlere uymayan klinik notların ve hastaların derlem içerisinde çıkartılması işleminden sonra derlem içerisinde 21,080 hasta ve bu hastalara ait 178,251 tane klinik not kalmıştır. Bu klinik notlar, Khadanga vd. [65] çalışmasında anlatılan ve kodu paylaşılan yöntemle veri ön işleme tabii tutulmuştur. Bu işlemler içerisinde, raporları bölümlere ayırma, cümle ve kelimeleri bölme ve token olarak ayırma gibi işlemler

mevcuttur. İlgili veri ön işleme koduna bu adresten ⁷ ulaşılabilir. Son olarak ise Bölüm 4.3’de anlatılan klinik varlık isim tanıma çalışması (med7 [53]) klinik notlar üzerine uygulanarak medikal terimler çıkartılmıştır. med7 çalışması tarafından eğitilen klinik alana özel VIT çalışması sayesinde, klinik notlar içerisindeki ilaç ismi, ilaç alım bilgisi, kullanım sıklığı, doz bilgisi, ilacın gücü, formu ve ilaç kullanımına devam etme süresi gibi medikal terimler çıkarılabilir olmuştur. Model eğitimleri sırasında veriler %70, %10, %20 oranlarında eğitim, validasyon, ve test kümelerine ayrılmışlardır. MIMIC-III ve deneyler içerisinde kullanılan derlemin istatistikleri Çizelge 5.1’de gösterilmiştir. Üzerinde çalışılan problemlerin kısa tanımları ve sınıf dağılım oranları da aşağıda okuyucular ile paylaşılmıştır.

- Hastanede mortalite: YBÜ ziyareti sonrası hastanede ölen hastalar (%10.5)
- YBÜ’de mortalite: YBÜ ziyareti esnasında YBÜ’de ölen hastalar (%7)
- YBÜ’de kalma süresi >3: YBÜ’de 3 günden fazla süre kalan hastalar (%43.2)
- YBÜ’de kalma süresi >7: YBÜ’de 7 günden fazla süre kalan hastalar (%7.9)

Zaman serisi öznitelikler ile model eğitimi. Hastaya ait yaşamsal gözlem verilerini ve laboratuvar sonuçlarını içeren 24 saatlik zamana bağlı öznitelikler ile LSTM ve GRU yöntemleri kullanılarak deneyler gerçekleştirilmiştir. Sonuçlar incelendiğinde, GRU mimarisi ile eğitilen modelin kesinlik (precision) değeri 0.56 ve duyarlılık (recall) değeri 0.34 olarak hesaplanmıştır. Daha açık bir şekilde ifade etmek gerekirse, sadece zaman serisi öznitelikleri ile eğitilen modellerin tahminlerinde, modelin mortalite olacak şekilde tespit ettiği hastaların %56’sı mortalite olmuştur. Doğruluk (accuracy) değeri ise 0.89 olup, yoğun bakımda kalan hastaların mortalite olacağı veya olmayacağı 0.89 doğruluk oranı ile tespit edilmiştir. Tez kapsamındaki çalışmalarda ise, hem üzerinde çalışılan veri setinin imbalans olması, hem de doğruluk, kesinlik, duyarlılık gibi metriklerin belirli eşik değerlerine göre hesaplanmasından kaynaklanan farklılıklarından ötürü model karşılaştırmaları AUROC ve AUPRC metrikleri üzerinden gerçekleştirilmiştir. Yapılan deneyler incelendiğinde, LSTM yöntemine göre daha basit bir mimari yapıya sahip GRU yönteminin AUROC ve AUPRC metrikleri bazında LSTM yöntemine göre %0.5-%1 daha yüksek performans gösterdiği gözlemlenmiştir. Bu sebeple çalışma boyunca yapılan bütün deneylerde, zamana bağlı öznitelikler GRU yöntemi kullanılarak işlenmiştir. Yalnızca zamana bağlı öznitelikler ile yapılan deneylerde, 256 nöronlu tek katmanlı bir GRU yapısı kullanılmıştır. Ayrıca önerilen GRU modelinin son

⁷https://github.com/kaggarwal/ClinicalNotesICU/blob/master/scripts/extract_notes.py

Çizelge 5.1: MIMIC-III veri setinin ve bu çalışmada kullanılan veri setinin istatistikleri.

	Hasta Sayısı	Hastane Bas, vuru Sayısı	YBÜ Bas, vuru Sayısı
MIMIC-III	46,520	58,976	61,532
MIMIC-III (> 15 yas, ından büyükler)	38,597	49,785	53,423
MIMIC-Extract	34,472	34,472	34,472
MIMIC-Extract (en az 24 saat YBÜ' de kalan hastalar)	23,937	23,937	23,937
Kullanılan Veri Seti (Medikal terimi olmayan hastaların elenmesinden sonra)	21,080	21,080	21,080

katmanında sigmoid aktivasyon fonksiyonu kullanılarak, modelin ikili sınıflandırma işlemlerini gerçekleştirilmesi sağlanmıştır.

Çok kipli Yaklaşım (Multimodal Approaches). Çok-kipli modellerin kullanıldığı deneylerde, zaman serisi özniteliklerinin yanısıra klinik notlar içerisinde çıkarılan medikal terim bilgileri de kullanılarak, klinik problemlerin çözümü için eğitilen modellerin performansı artırılmaya çalışılmıştır. Bu çalışmada kapsamında, önceki çalışmamızda [70] olduğu gibi doğrudan klinik notları model içerisinde kullanmak yerine, klinik notlar içerisinde medikal terimleri çıkartarak, bu terimlerin model içerisinde kullanılması hedeflenmiştir. Literatürdeki klinik alana özgü varlık isim tanıma modelleri [52, 53, 125] arasından MIMIC-III veri seti ile eğitilmiş olan med7 [53] çalışması, medikal terimlerin klinik notlar içerisinde çıkarılabilmesi için seçilmiştir. Med7 çalışması tarafından önerilen ve paylaşılmış olan klinik varlık isim tanıma yöntemi, MIMIC-III içerisindeki klinik notlara uygulandığında çıkarılan medikal terimler, bu terimlerin sayısı ve örnekleri Çizelge 5.2'de paylaşılmıştır.

Kelime Temsili. Klinik notlar içerisinde çıkarılan medikal terimlerin derin öğrenme tabanlı modellere girdi olarak verilebilmesi için kelime temsil yöntemleri kullanılarak vektörel forma dönüştürülmesi gerekmektedir. Farklı kelime temsil yöntemleri aynı kelime için farklı semantik öznitelikleri yakalayabilmektedir. Bu sebeple, deneylerde detayları Bölüm 4.2'de anlatılan Word2Vec [19], FastText [36] ve bu iki temsil yönteminin ürettiği vektörlerin birleştirilmiş hali kullanılmıştır. Çalışmamız kapsamında kullanılan klinik alana özgü Word2Vec ve FastText modelleri, Huang vd. [48] tarafından yapılan çalışmada önerilmiştir. Bu modellerin eğitiminde yaklaşık 2.8 milyar kelimedenden oluşan MIMIC-III klinik not derlemi kullanılmıştır. Her iki temsil yöntemi de kelimelerin 100 boyutlu ($w_i \in \mathbb{R}^{100}$) vektörel temsillerini üretmektedir. Son olarak ise hem Word2Vec hem de FastText yöntemlerinden yararlanabilmek adına, iki yöntemin vektörel çıktıları yatay olarak birleştirilerek ($c_i \in \mathbb{R}^{200}$) üçüncü bir yöntem önerilmiştir. Eğitilmiş Word2Vec modelinde ilgili medikal kelimenin bulunamaması durumunda, tamamı sıfırlardan oluşan (zero padding) 100 boyutlu vektörel temsil kullanılmıştır. **Doküman Temsili.** Kelime temsil yöntemlerine ek olarak, doküman temsil etme yöntemi olan Doc2Vec [68] (detaylar için Bölüm 4.2.3 incelenebilir.) yönetimi ile

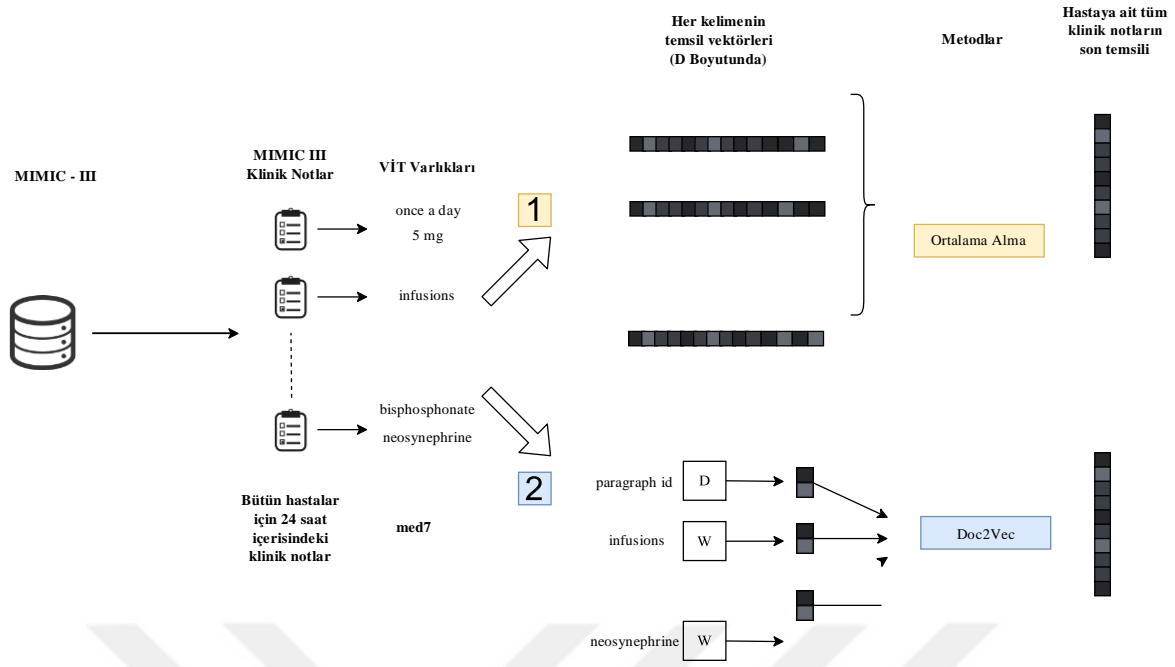
Çizelge 5.2: İlk sütun medikal terimlerin tiplerini belirtirken, ikinci sütun derlem içerisindeki klinik notlarda bu terimlerin isimlerinin kaçar kez geçtiği göstermektedir. Üçüncü sütun, bu ifadelerin tekil olarak kaçar kez gösterirken, son sütunda ise her bir terim için örnek ifade verilmiştir.

Varlık İsmi	Toplam Sayı	Tekil Sayı	Örnek Terim
İlaç ismi	744778	18268	Tacrolimus
İlaç gücü	156486	10749	400mg/5ml
İlaç formu	40885	597	suspension
İlaç alım bilgisi	207876	1193	PO (twice a day)
İlaç doz bilgisi	126756	7239	30ml
İlaç kullanım sıklığı	71285	3344	bid
İlaç devam süresi	5939	1185	next 5 days

de deneyler gerçekleştirilmiştir. Bir hasta için çıkartılan medikal terimlerin tamamı bir doküman olarak varsayılmış ve Doc2Vec algoritmasına girdi olarak verilmiştir. Bağlam pencere uzunluğu 5 kelime olarak seçilen Doc2Vec yöntemi ile, her hasta için hastaya ait medikal terimlerin tamamı dikkate alınarak sabit uzunlukta 100 boyutlu bir vektör üretilmiştir. Hem kelime temsil yöntemleri kullanılarak öğrenilen hem de doküman temsil etme yöntemi kullanılarak öğrenilen temsiller Şekil 5.1’de paylaşılmıştır.

Bu çalışmada, medikal terimler üzerinden üç farklı şekilde öznitelik vektörü çıkartılarak farklı çok-kipli model eğitimi yapılmıştır. Medikal terimlerin öznitelik vektörü haline dönüşümünde, kelime temsillerinin ortalaması, doküman temsil yöntemi (Doc2Vec) ve çalışmada kapsamında en iyi sonucu veren 1D CNN yöntemi kullanılmıştır.

Kelime Temsillerinin Ortalaması Yöntemi ile Çok-kipli Modeller. Önerilen bu model, hastalara ait klinik notlar içerisinde çıkartılan medikal terimlerin vektörel ortalamasını girdi alacak şekilde kurgulanmıştır. Her bir hastanın N adet klinik notu ve bu klinik notlar içerisinde çıkartılan K adet medikal terim bulunabilmektedir. Her bir medikal terim Bölüm 4.2’de anlatılan kelime temsil yöntemleri ile temsil edilmiştir. Hastaya ait klinik notlar içerisinde çıkartılan K adet n boyutlu vektör

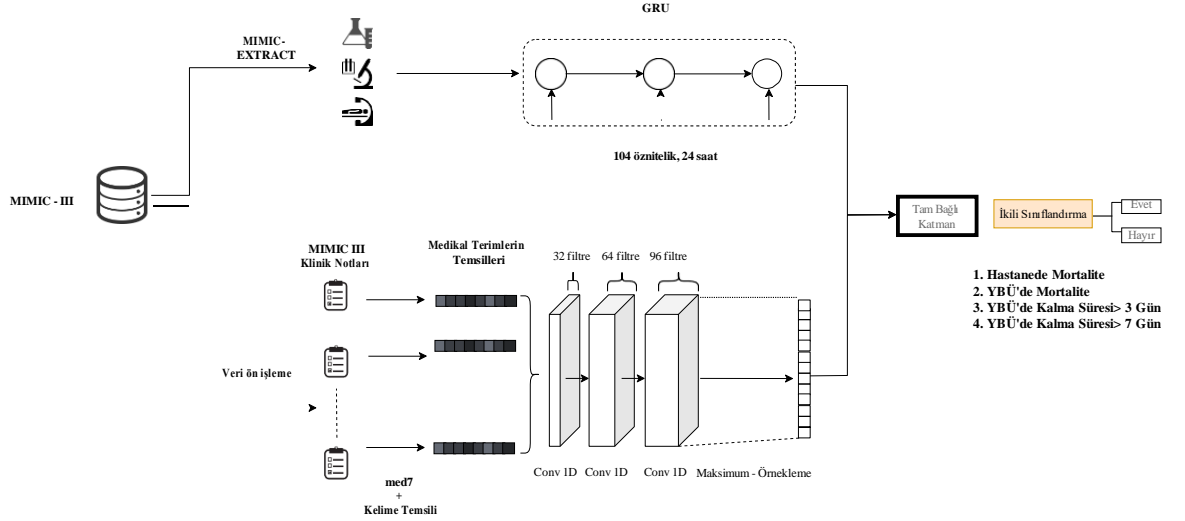


Şekil 5.1: Medikal varlık isim vektörlerini öğrenmek için önerilen yöntem. (1) Klinik notlar içerisinde çıkarılan medikal varlıklar sürekli değerler içeren vektörlere dönüştürülür. Daha sonra ise öğrenilmiş temsillerin ortalaması alınmıştır. (2) Eğer klinik notlar içerisindeki kelimelerden biri medikal varlık tiplerinden herhangi birine girmiyorsa o kelimeler klinik notlar içerisinde silinmiştir. Ardından, veri ön işleme uygulanmış, klinik notlar üzerinden Doc2Vec yöntemi eğitilmiştir.

öncelikle toplanmış, daha sonra ise K sayısına bölünmüştür. Bu sayede medikal terim temsillerinin ortalaması alınmıştır. Zamana bağımlılık, 256 nöron içeren GRU ile işlenmiş, ve çıktısı kelime temsillerinin ortalaması ile birleştirilerek 256 nörondan oluşan yapay sinir ağına girdi olarak verilmiştir.

Doc2Vec Yöntemi ile Çok-kipli Modeller. Önerilen bu yöntemde, kelime temsillerinin ortalaması almak yerine, Bölüm 4.2.3’de detaylandırılan Doc2Vec yöntemi kullanılarak, K adet medikal terim için sabit-uzunlukta öznitelik vektörü öğrenilmiştir. Bunun için, öncelikle hastaya ait N adet klinik not birleştirilmiş, ve içerisinde medikal terim olmayan kelimeler önceden eğitilmiş klinik varlık isim tanıma algoritmasıyla çıkarılmıştır. Bu işlemden sonra ise hastaya ait medikal terimlere Doc2Vec algoritması uygulanarak, her hasta için klinik notlarının temsilleri öğrenilmiştir. Doc2Vec yöntemi ile sabit-uzunlukta öğrenilen öznitelik vektörlerinden sonra ise önerilen çok-kipli derin öğrenme mimarisi, "kelime temsillerinin ortalaması" yöntemi ile beraber önerilen çok-kipli derin öğrenme mimarisinin aynı olmuştur. Her iki yöntemin özet çizimi Şekil 5.1’de okuyucu ile paylaşılmıştır.

Önerilen Çok-kipli Modeller. Çalışma kapsamında en iyi sonucu veren ve önerilen yöntem Şekil 5.2’de okuyucu ile paylaşılmıştır. Bu yöntemde, çok-kipli model, medikal kelimelerin temsilleri üzerinden öznitelik çıkarılması için 1D CNN yönteminden



Şekil 5.2: Hastane içi mortalite, YBÜ mortalite, Hastanede Kalma Süresi >3, ve Hastanede Kalma Süresi >7 klinik problemleri için önerilen yöntem özetleri. MIMIC-III içerisinde zamana bağlı özelliklerin hazırlanması için MIMIC-Extract çalışması kullanılmış, ve özellikler GRU algoritmasına girdi olarak verilmiştir. Aynı zamanda klinik notlara veri ön işleme adımları uygulanmış, ve medikal terimlerin çıkartılması için med7 yöntemi kullanılmıştır. Medikal terim temsilleri içerisinde özellik çıkartımı için 1D CNN yöntemi önerilmiştir. Son olarak ise 4 farklı sınıflandırma problemini tahminleyebilmek için çıkartılan özellikler birleştirilmiş ve 3 katmanlı yapay sinir ağına girdi olarak verilmiştir.

faydalanmaktadır. Detayları Bölüm 4.1.4’de anlatılan 1D CNN yöntemi, metinsel veriye uygulandığında yan yana geçen komşu kelimeler arasındaki ilişkiyi de öğrenerek doğal dil işleme alanındaki problemlerde başarılı sonuçlar vermektedir [38]. Önerilen yöntemde, K tane medikal kelime, hastaya ait N tane klinik not içerisinde çıkartılmıştır. Çıkartılan K adet medikal kelime, öncelikli olarak Word2Vec, FastText ve bu iki yöntemin birleştirilmesi ile beraber üç farklı şekilde temsil edilmiştir. Daha sonra ise bu medikal kelime temsilleri $e_i \in \mathbb{R}^d$, dikey olarak birleştirilerek her hasta için bir matris $M \in \mathbb{R}^{k \times d}$ oluşturulmuştur. Gerekli olduğu zamanlarda sıfırlardan oluşan vektörlerle tamamlanan bu klinik varlık isim tanıma matrisi s_u şeklinde temsil edilmiştir:

$$e_{1:k} = e_1 \otimes e_2 \otimes \dots \otimes e_k \quad (5.1)$$

\otimes sembolü, vektörleri dikey biçimde birleştirme operasyonunu temsil ederken, e sembolü her bir medikal terimi temsil etmektedir. k ise medikal terim sayısını belirten semboldür. Yapılan bu çalışmada, Öztürk vd. [57] tarafından yapılan çalışmada kullanılan 1D CNN modeline benzer bir model kullanarak medikal terimler içerisinde özellikler çıkartılmıştır. 32, 64, ve 96 filtre büyüklüğüne sahip 3 adet 1D evrimsel

katmanlar arka arkaya koyarak oluşturulan bu mimarinin son katmanında maksimum-örnekleme (max-pooling) yöntemi kullanılmıştır. Maksimum-örnekleme ile çıkartılan öz nitelikler, GRU tarafından çıkartılan öz nitelikler ile birleştirilerek 512 nörondan oluşan yapay sinir ağına girdi olarak verilmiştir.

5.3 Deneysel Sonuçlar

Bu bölümde yapılan deneylerde elde edilen sonuçlar raporlanmıştır. Ayrıca, modellerin değerlendirilmesinde kullanılan metrikler ve deney detaylarından da bahsedilmiştir. Deney Ayarları. Üzerinde çalışılan bütün klinik problemlerin model eğitimlerinde, hastanın YBÜ'deki ilk 24 saatteki verileri kullanılmıştır. Çok-kipli modellerde, model mimarisinin yapay sinir ağı katmanında 0.2 oranında seyreltme (dropout) [126] kullanılmıştır. Modellerin doğrusal olmayan ilişkileri öğrenebilmesi adına aktivasyon fonksiyonu olarak Rectified Linear Unit (ReLU) [127] yöntemi kullanılmıştır. Aşırı öğrenmeyi engellemek için ise regülarizasyon yöntemlerinden L_2 kullanılmış, ve regülarizasyon ölçeği 0.01 olarak seçilmiştir. Hata fonksiyonunun optimizasyonu için Adam [128] yöntemi 0.001 öğrenme oranıyla tercih edilmiştir. Bütün modeller ikili çapraz-entropi (binary cross-entropy) ve bağımsız parametre optimizasyonu ile eğitilmiştir. Hiper parametre optimizasyonu sırasında ise, gizli katman sayıları, nöron sayıları, evrisimsel filtreler, filtre büyüklükleri, öğrenme katsayısı gibi birçok değişken üzerinde çalışılmıştır. Her model 50 iterasyon boyunca eğitilmiş ve erken durdurma (early stopping) yöntemi validasyon hatası (validation loss) üzerinde gözlemlenmiştir. Her problem için aynı modeller 10 farklı başlatma değeri (initialization seed) ile denenmiş ve ortalama değerler okuyucu ile paylaşılmıştır. Modelleri değerlendirmek adına Bölüm 4.6'de anlatılan AUROC, AUPRC ve F1 metrikleri kullanılmıştır. Modellerin eğitimi esnasında derin öğrenme yöntemleri için arka planda Tensorflow [122] ve Tensorflow kütüphanesini kullanan Keras [123] kütüphaneleri kullanılmıştır. Bütün deneyler NVIDIA Tesla K80 GPU, 24 GB VRAM, 378 GB RAM ve Intel Xeon E5 2683 işlemci içeren bilgisayar üzerinde gerçekleştirilmiştir. Çalışmaya ait kodlara <https://github.com/tanlab/ConvolutionMedicalNer> adresinden erişilebilir.

Temel Yöntemlerin (Baseline) Deney Sonuçları.

Temel yöntemler hem zaman serisi öz niteliklerini kullanılarak GRU eğitimini hem de Doc2Vec ve ortalama alma yöntemleri ile çok-kipli model eğitimlerini kapsamaktadır. Çizelge 5.3 bütün temel yöntemlerin sonuçları özetlemektedir. Sonuçlar incelendiğinde, sadece zaman-serisi öz niteliklerini kullanarak eğitilen GRU yönteminin başarılı sonuçlar vermesine rağmen, beklenildiği gibi çok-kipli yaklaşımın gerisinde kaldığı gözlemlenmektedir. Zaman serisi öz niteliklerine ek olarak medikal terimlerin kullanılması ile önerilen çok-kipli derin öğrenme modellerinin, hastane içi mortalite tahmininde, %1.5

Çizelge 5.3: Temel yöntemlerin performans karşılaştırması. 4 ayrı klinik problemin her biri için, AUC, AUPRC ve F1 skorlarının ortalamasına standart sapma değerleri raporlanmıştır.

Problem	Yöntem	Temsil Yöntemi	AUROC	AUPRC	F1
Hastanede Mortalite	GRU	-	85.04±0.004	52.15±0.009	42.29±0.016
	Doc2Vec ile Çok-kipli model	Doc2Vec	85.96±0.002	54.17±0.00446.60±0.016	
		Word2Vec	86.42±0.004	54.22±0.008	45.42±0.013
	Kelime Temsil Ortalaması ile Çok-kipli model	FastText	86.09±0.00454.47±0.007	45.50±0.010	
		Word2Vec+FastText	85.98±0.002	54.19±0.008	45.66±0.021
YBÜ'de Mortalite	GRU	-	86.32±0.004	46.51±0.011	36.30±0.026
	Doc2Vec ile Çok-kipli model	Doc2Vec	86.80±0.002	48.22±0.006	41.95±0.017
		Word2Vec	87.17±0.00248.47±0.006	42.30±0.021	
	Kelime Temsil Ortalaması ile Çok-kipli model	FastText	87.14±0.003	48.36±0.00642.91±0.014	
		Word2Vec+FastText	86.90±0.004	48.28±0.007	40.76±0.022
YBÜ'de Kalma Süresi > 3 Gün	GRU	-	67.40±0.003	60.17±0.005	53.36±0.016
	Doc2Vec ile Çok-kipli model	Doc2Vec	68.90±0.00261.88±0.002	54.32±0.008	
		Word2Vec	68.63±0.003	61.81±0.003	54.19±0.012
	Kelime Temsil Ortalaması ile Çok-kipli model	FastText	68.55±0.003	61.59±0.003	54.46±0.012
		Word2Vec+FastText	68.61±0.003	61.69±0.00354.70±0.009	
YBÜ'de Kalma Süresi > 7 Gün	GRU	-	70.54±0.004	16.35±0.0062.33±0.012	
	Doc2Vec ile Çok-kipli model	Doc2Vec	71.63±0.005	17.22±0.004	1.50±0.007
		Word2Vec	71.59±0.00517.91±0.006	1.35±0.008	
	Kelime Temsil Ortalaması ile Çok-kipli model	FastText	71.31±0.008	17.57±0.007	1.02±0.008
		Word2Vec+FastText	71.59±0.007	17.67±0.007	1.37±0.013

AUROC, %2.5 AUPRC ve %4 F1 skoru iyileşme sağladığı görülmektedir. Bir diğer mortalite tahmin problemi olan YBÜ'de mortalitede de AUROC ve AUPRC metriklerinde %2'lik bir iyileşme, F1 metriğinde ise %7'lik bir iyileşme mevcuttur. Mortalite problemlerinin yanısıra çok-kipli yaklaşımlar YBÜ'de kalma süresi tahmininde de model başarımlarını olumlu yönde etkilemiştir. Hem YBÜ'de > 3 günden fazla kalma,

hem de YBÜ’de > 7 günden fazla kalma problemlerinde bütün metrikler yaklaşık olarak %1.5 oranında artış göstermiştir.

Üzerinde çalışılan klinik problemlerin tahminini iyileştirmek için medikal terimleri kullanmanın avantajını göstermek adına ortalama alma ve Doc2Vec yöntemleri ile çok-kipli modeller eğitilmiştir. Alınan sonuçlar ise açık bir şekilde, medikal terimleri zaman-serisi öznitelikleri ile beraber kullanmanın sonuçları iyileştirdiğini göstermektedir. Doc2Vec yöntemi ile ortalama alma yöntemleri karşılaştırıldığında ise aralarında önemli bir başarı farkının olmadığı gözlemlenmektedir. Bu sebeple medikal terimler üzerinden daha etkili öznitelik çıkartabilmeyi hedefleyen evrimsel tabanlıların öğrenme mimarisi önerilmiştir.

Önerilen Yöntemin Deney Sonuçları. Önerilen evrimsel tabanlı çok-kipli derin öğrenme modelinin sonuçları, çalışma kapsamında yapılan diğer deneyler ile karşılaştırılmış ve bu sayede önerilen modelin etkinliği ve güvenilirliği tartışılabilmiştir. Çizelge 5.4’de paylaşılan sonuçlar incelendiğinde, önerilen evrimsel tabanlı çok-kipli model sonuçlarının diğer çok-kipli modellere göre daha iyi başarı gösterdiği görülmektedir. YBÜ’de 7 günden fazla kalma problemindeki F1 skoru haricinde diğer bütün problemlerde ve metriklerde önerilen çok-kipli modelin en iyi sonucu verdiği ve diğer bütün modellerden daha başarılı olduğu görülmektedir. Sonuçlar incelendiğinde bütün problemlerde ve bütün metrikler için, ortalama %1-%2 başarı iyileşmesi sağlanarak, evrimsel katmanların, ortalama alma ve Doc2Vec yönteminden daha iyi öznitelik çıkarımını gerçekleştirdiği söylenebilmektedir.

5.4 Değerlendirme

Çizelge 5.3’de görüldüğü üzere çok-kipli yaklaşım bütün klinik problemlerde başarıyı olumlu yönde etkilemiştir. Yapılan deneyler ayrıca, medikal terimlerin hangi kelime temsil etme yöntemi ile temsil edilmesi gerekliliği için de önemli bir fırsat sağlamaktadır. Sonuçlar incelendiğinde ise, kelime temsilleri içerisinde kesin bir kazanan bulunmadığı farklı problemler ve farklı metriklere göre değişkenlik gösterdiği gözükmektedir. Buna rağmen genel olarak Word2Vec modelinin en iyi sonuç veren yöntem olduğu (genellikle %0.5 kadar) söylenebilmektedir. Kelime temsillerinin karşılaştırılmasında yöntem karşılaştırmasında da, ortalama alma yöntemi ile Doc2Vec yöntemi karşılaştırıldığında, ortalama alma yönteminin genel olarak daha iyi sonuç verdiği gözlenmektedir. Bu karşılaştırmaların yanı sıra, çalışmanın temel amaçlarından bir tanesi medikal kelimelerin efektif bir şekilde özniteliklere dönüştürülmesidir. Ortalama alma ve Doc2Vec yöntemleri, medikal terimleri kullanmanın avantajını ortaya çıkartmasına rağmen, medikal terimlerin vektörleri üzerinden daha efektif bir şekilde öznitelik çıkartabilmek adına 1-boyutlu evrimsel sinir ağlarından yararlanılmıştır. Üç adet evrimsel bloku peş peşe

Çizelge 5.4: Önerilen yöntem ile diğer en iyi temel yöntemlerin karşılaştırılması.

Problem	Yöntem	Temsil Yöntemi	AUROC	AUPRC	F1
Hastanede Mortalite	En iyi temel yöntem (best baseline)	-	86.42±0.004	54.47±0.007	46.60±0.016
		Word2Vec	87.55±0.003	55.87±0.008	47.23±0.014
	Önerilen Model	FastText	87.15±0.002	55.68±0.005	46.87±0.015
		Word2Vec+FastText	86.98±0.003	55.35±0.008	46.38±0.027
YBÜ'de Mortalite	En iyi temel yöntem	-	87.17±0.002	48.47±0.006	42.91±0.014
		Word2Vec	88.35±0.002	49.23±0.008	43.02±0.029
	Önerilen Model	FastText	87.85±0.001	48.78±0.009	43.09±0.026
		Word2Vec+FastText	87.66±0.002	48.74±0.009	42.24±0.027
YBÜ'de Kalma Süresi>3 Gün	En iyi temel yöntem	-	68.90±0.002	61.88±0.002	54.70±0.009
		Word2Vec	69.54±0.002	62.68±0.003	55.04±0.012
	Önerilen Model	FastText	69.61±0.003	62.55±0.003	55.87±0.017
		Word2Vec+FastText	69.93±0.001	62.77±0.002	55.82±0.008
YBÜ'de Kalma Süresi>7 Gün	En iyi temel yöntem	-	71.63±0.005	17.91±0.006	1.33±0.012
		Word2Vec	72.55±0.005	18.78±0.006	1.58±0.001
	Önerilen Model	FastText	71.81±0.004	18.01±0.004	1.08±0.008
		Word2Vec+FastText	71.92±0.007	18.25±0.006	1.38±0.009

ekleyerek oluşturulan bu yapının sonuna 1-boyutlu bir maksimum-seyreltme operasyonu uygulanarak sabit uzunlukta vektör elde edilmesi sağlanmıştır. Çizelge 5.4'de paylaşılan sonuçlar incelendiğinde ise önerilen yöntemin diğer çok-kipli modellerden de başarılı sonuç verdiği görülmektedir.

Literatürde, mortalite ve yoğun bakımda kalma süresini tahmin etmeye yönelik birçok çalışma bulunmaktadır. Purushotham vd. [85] çalışmalarında hastane içi mortalite, yoğun bakımda kalma süresi ve ICD-9 kod grupları tahmini problemleri üzerine çalışmıştır. MIMIC-III veri seti içerisinde, hastalara ait 12 klinik öznetelik çıkartmak

için bir veri ön işleme yöntemi önermektedir. Ayrıca çıkartılan bu öznitelikler, geleneksel yapay öğrenme teknikleri ile birleştirilerek çeşitli klinik problemler için tahmin işlemi gerçekleştirilmektedir. Bu çalışmada ise daha güncel ve daha büyük öznitelik kümesi çıkartan MIMIC-Extract çalışması kullanılmıştır. Ayrıca, çalışmalarımızda üzerinde çalışılan problemlerin tamamını sınıflandırma problemi olarak tasarlanırken, Purushotham vd.'nin yaptığı çalışmada YBÜ'de kalma problemi regresyon problemi olarak tasarlanmıştır. Bir diğer çalışmada Jin vd. [67] çok-kipli bir derin öğrenme yöntemi önererek zaman-serisi öznitelikler ile medikal terimleri birleştirilerek hastane içerisindeki mortaliteyi tahmin etmeye çalışmaktadır. Klinik notları temsil etmek için Document Vector through Corruption (Doc2VecC) [68] yöntemi önerilmiştir. Bu çalışmada ise, sadece hastane içi mortalite tahmini yerine YBÜ'de mortalite ve YBÜ'de kalma sürelerini de tahmin edilmiştir. Buna ek olarak, medikal terimlerin temsili için hem Word2Vec, FastText gibi farklı kelime temsil yöntemleri hem de bu temsiller üzerinden öznitelik çıkartabilecek ortalama alma, Doc2Vec, evrimsel sinir ağları gibi yöntemler denenerek kapsamlı bir çalışma gerçekleştirilmiştir.

Sonuç olarak, önerilen yöntem sayesinde, diğer çok-kipli yöntemlerden AUPRC metriği bazında ortalama %1 -%1.5, zaman-serisi modellerine göre ise AUPRC metriği bazında %2.5-%3 iyileşme görülmüştür. Bu iyileştirmelere rağmen çalışma kapsamında çeşitli kısıtlar bulunmaktadır. İlk olarak, kelime temsil yöntemi olarak sadece bağlam bağımsız Word2Vec ve FastText gibi yöntemler kullanılmıştır. İkinci olarak, çalışmada kullanılan kelime temsil yöntemleri ve varlık isim tanıma modeli MIMIC-III veri seti üzerinden eğitilmiştir. Bu durum, önerilen yöntemin bir başka ESK veri seti üzerinde doğrudan çalışmasını kısıtlayabilecek bir nokta olmaktadır. Son olarak ise, önerilen yöntemin açıklanabilirliği kısıtlıdır. Yapılan tahminlerin neye dayanarak yapıldığının açıklanamaması, modeli kullanacak olan sağlık personeli için kritik bir konu olarak değerlendirilmektedir.

Bu sebeple bu çalışma farklı şekilde genişletilebilir. Öncelikle, bağlam-bağımlı BERT ve benzeri yöntemler kullanılarak temsil edilen medikal terimlerin başarı artırma ihtimali öngörülmektedir. İkinci olarak, sadece MIMIC-III üzerinden değil, daha geniş bir derlem üzerinden eğitilen klinik kelime temsilleri ve klinik varlık isim tanıma yöntemi ile daha güçlü ve başarılı sonuçlar alınabileceği düşünülmektedir. Hastaya ait ilaç verileri, radyolojik görüntüler gibi bilgilerinde modele dahil edilmesi ile beraber tahmin başarımının artabileceği öngörülmektedir. Son olarak ise modellerin açıklanabilirliğinin artması için modele dikkat mekanizması eklenmesi veya literatürde önerilen diğer açıklanabilirlik yöntemleri ile (örn: SHapley Additive exPlanations (SHAP) [75]) açıklanabilir tahminler elde edilmesi gelecekte gerçekleştirilebilecek konular olarak görülmektedir.

6. ZAMAN SERİSİ ve KLİNİK NOTLAR ile TAHMİN ETME

6.1 Motivasyon

Literatürde, hastaya ait yaşamsal gözlem verileri ve laboratuvar sonuçları gibi zamana bağlı öznitelikler ile hastanın mortalite olması ve yoğun bakımda ne kadar süre kalacağı tahmin eden çeşitli çalışmalar mevcuttur. Bu çalışmada ise hastaya ait zamana bağlı özniteliklere ek olarak, hastaya ait kıymetli bilgiler içeren klinik not ve raporlarında kullanılması önerilmiştir. Bölüm 4.2’de detayları anlatılan Word2Vec ve FastText gibi bağlam bağımsız (context-free) kelime temsillerinin ortaya çıkması ile beraber birçok alanda metinsel veriler çok daha efektif bir şekilde kullanılmaya başlanmıştır. Bu çalışmalara ek olarak, doğal dil işleme alanında çığır açan, geleneksel kelime temsil yöntemlerinin dezavantajlarını kapatan ve hemen hemen bütün doğal dil işleme probleminde literatürdeki en iyi sonucun alınmasını sağlayan dönüştürücü (transformer) tabanlı BERT modeli ile beraber metinsel veriler üzerine yapılan çalışmalar sayısının artmasıdır.

Üzerinde çalışılan, mortalite ve YBÜ’de kalma süresi problemleri, dört farklı sınıflandırma problemi olarak ele alınmıştır. Bu problemler, hastanede mortalite, yoğun bakım ünitesinde (YBÜ) mortalite, YBÜ’de 3 günden fazla kalma ve YBÜ’de 7 günden fazla kalma olarak ele alınmıştır. Önerilen derin öğrenme tabanlı çok-kipli model ile, hastaya ait ilk 24 saatlik yapısal zamana bağlı özniteliklerle beraber yapısal olmayan klinik notları aynı anda kullanarak üzerinde çalışılan klinik problemlerde başarılı sonuçlar alınması hedeflenmiştir. Alınan sonuçlar incelendiğinde klinik notları kullanmanın, model başarımlarının nasıl değiştiği ortaya konmuştur.

6.2 Önerilen Yöntem

Hastaya ait zamana bağlı öznitelikler Bölüm 3.2’de anlatılan MIMIC-Extract çalışması ile elde edilmiştir. Çalışma kapsamındaki eğitilen ilk modelde, hastaya ait 104 (detaylar için Bölüm 3.2 incelenebilir.) adet yaşamsal gözlem ve laboratuvar sonuç öznitelikleri kullanılmıştır. Bu 104 öznitelik, hastaya ait ilk 24 saatlik süre zarfında toplanmış ve Bölüm 4.1.3’de detayları anlatılan Geçitli tekrarlayan birim (Gated Recurrent Unit, GRU) yapısı kullanılarak eğitilmiştir. Eğitilen GRU modeli tek katmanlı ve 128 nöron

çerçek s,ekilde tasarlanmıs,tır. GRU modelinin sonuna sigmoid fonksiyonu eklenerek ikili sınıflandırma is,lemi gerçekles,tirilmiştir.

Çalışmalarda kullanılan bir diđer veri türü ise klinik notlardır. MIMIC-III içerisinde (NOTEEVENTS.csv) 15 farklı klinik not çeş,idi bulunmaktadır. MIMIC-III veri seti içerisinde toplamda 46,520 hastaya ait 2,083,180 klinik notu bulunurken bu klinik notların kategorileri ve her kategoriye ait not sayısı Çizelge 6.1’de okuyucu ile paylas,ılmış,tır. MIMIC-Extract veri ön is,lemesi sonrasında çeş,itli veri ön is,leme ve hasta eleme kriterleri sebebi ile derlem içerisinde çıkartılan hastalar olduđu gibi, klinik notlar içinde üç ayrı kriter belirlenerek çeş,itli klinik notlar ve klinik notlara sahip olmayan hastalar derlem içerisinde çıkartılmış,tır. Bu kriterler aşağıda listelenmiştir:

- Hasta taburcu notu (Discharge summary) içerisinde mortalite veya yoğun bakımda kalma süresine ait doğrudan bir bilgi içerebileceđi için, bu notların tamamı çıkartılarak olası veri kaça ğı sorununu engellenmeye çalışılmış,tır.
- Zamana bađlı özniteliklerin seçimindeki zaman kısıtlı klinik notların seçiminde de kullanılmış,tır. Bu sebeple, hasta yoğun bakıma yattıktan sonraki ilk 24 saatte, hastaya yazılan klinik notlar seçilmiş, ve sadece bu notlar modele girdi olarak verilmiş,tır.
- Bu iki kriterin veri setine uygulanmasından sonra ilk 24 saat içerisinde hiç klinik notu bulunmayan hastalar derlem içerisinde çıkartılarak deneylere dahil edilmemiş,tır.

46,520 hasta verisi içeren MIMIC-III veri setine MIMIC-Extract veri ön is,leme çalış,ması uygulandıđında hasta sayısı 23,937’e düşmektedir. Ardından, ilk 24 saat içerisinde klinik notu olmayan hastalarda derlem içerisinde çıkartıldıđında deneylerde verisi kullanılan hasta sayısı 21,087’ye düşmüştür. Her adımdan sonra oluşan hasta sayısı ve diđer bilgiler özet halinde Çizelge 6.2’de okuyucu ile paylas,ılmış,tır. Ayrıca, mevcut 2,083,180 klinik not sayısı, hasta elenmesine bađlı olarak 178,251’e düşmüş ve nihai modellerde bu klinik notlar kullanılmış,tır.

Klinik notlar doğası gere ği karmaşık, içerisinde jargon ve kısaltmalar içeren metinlerdir. Bu sebeple klinik notlar, vektörel formlara dönüş,türülmeden önce veri ön is,leme adımlarından geçirilmiş,tır. Bu veri ön is,leme adımları, bos,lukları temizleme, "doktor, dr, admission date:" vb. ifadeleri silme gibi is,lemler olmakla beraber çeş,itli kısaltmalar da uzun haline dönüş,türülmüş,tür. Örnek olması açısından, notlar içerisinde geçen çeş,itli kısaltmalar ve bu kısaltmaların hangi kelimelere dönüş,türüldüđu Çizelge 6.3’de paylas,ılmış,tır.

Veri ön is,leme is,leminden sonra klinik notların vektörel hale dönüş,türülmesi için çalış,ma yapılmıştır. Klinik notların temsili, önceden eğitilmiş ve detayları Bölüm 4.2.4’de

Çizelge 6.1: Önerilen model performansının, temel olarak alınan modellerle karşılaştırılması. Her problem ve metrik için en yüksek skorlar vurgulanmıştır.

Kategori İsmi	Kategori İsmi (Orijinal)	Not Sayısı
Hemşire/Diğer	Nursing/other	822,497
Radyoloji	Radiology	522,279
Hemşire	Nursing	223,556
Elektrokardiyogram	ECG	209,051
Doktor	Physician	141,624
Hasta taburcu notu	Discharge Summary	59,652
Ekoensefalogram	Echo	45,794
Solunum	Respiratory	31,739
Beslenme	Nutrition	9,418
Genel	General	8,301
Rehabilitasyon Hizmetleri	Rehap Services	5,431
Sosyal Hizmet	Social Work	2,670
Vaka Yönetimi	Case Management	967
Eczane	Pharmacy	103
Danışman	Consult	98

anlatılmış, olan ClinicalBERT modeli ve Sentence-Bert (Detaylar için Bölüm 4.2.6) yöntemi kullanılarak gerçekleştirilmiştir. Her bir klinik not 768 boyutlu vektörler ile temsil edilmiştir. Bir hastaya ait birden fazla klinik not olabileceği için, hastaya ait klinik not temsillerinin ortalaması alınarak nihai temsil ortaya çıkartılmıştır. Zaman serisi verileri ile elde edilen 128 boyutlu temsiller ile klinik notlar ile elde edilen 768 boyutlu temsiller birleştirilerek sırasıyla 1024 ve 512 nörondan oluşan 2 katmanlı bir yapay sinir ağına girdi olarak verilmiştir. Önerilen derin öğrenme tabanlı çok-kipli

Çizelge 6.2: MIMIC-III ve bu çalışmada kullanılan veri seti istatistikleri.

	Hasta Sayısı	Hastane Bas,vuru Sayısı	YBÜ Bas,vuru Sayısı
MIMIC-III	46,520	58,976	61,532
MIMIC-III (>15 yaş,ından büyükler)	38,597	49,785	53,423
MIMIC-Extract	34,472	34,472	34,472
MIMIC-Extract (en az 24 saat YBÜ kalan hastalar)	23,937	23,937	23,937
Kullanılan Veri Seti(Klinik notu olmayan hastaların elenmesinden sonra)	21,087	21,087	21,087

modelin detaylıtasarımıŞekil 6.1'de sunulmaktadır.

Yapay sinir ağıkatanlarında, aktivasyon fonksiyonu olarak Rectified Linear Unit (ReLU) [129] kullanılmış,tır. Model optimizasyonu için Adam [128] yöntemi seçilmiş, ve öğrenme oranı0.01 olarak belirlenmiştir. Her klinik problem deneyi için, modeller 100 iterasyon (epoch) boyunca eğitilmiştir. Aşırıö ğrenme (overfitting) probleminin önüne geçebilmek adına ise validasyon hatası takip edilerek erken durdurma (early stopping) yöntemi uygulanmıştır. Model başarımlarıdi ğer çalışmalarda da olduğu ve Bölüm 4.6'da detaylarıaçıklanmış,üzere AUROC, AUPRC ve F1 metrikleri üzerinden gerçekleştirilmiştir. Modellerin başlangıç ağırlıklarında veya diğ er çeşitli noktalarda oluş,abilecek rastgelelilikten kaynaklanabilecek farklarıönleyebilmek adına bütün deneyler 10 kez eğitilmiş ve ortalama sonuçlar paylaşılmıştır.

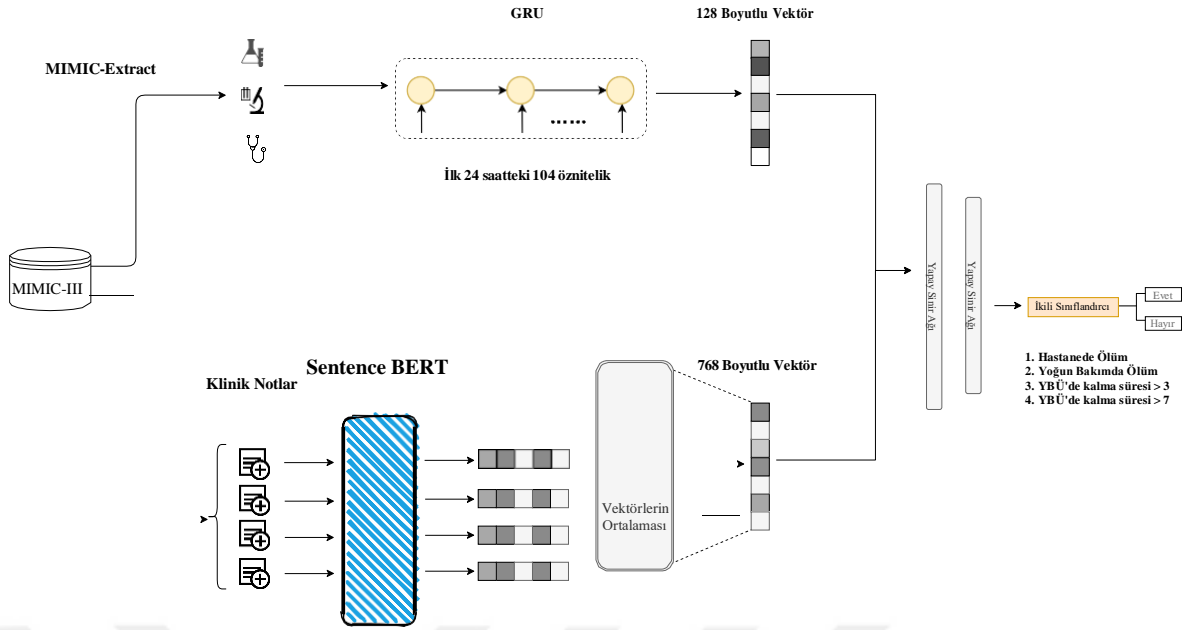
Uygulama Detayları.Bu çalışma kapsamında e ğitilen derin öğrenme tabanlımodeller, Keras [130] kütüphanesi ile gelis,tirilmiştir,tır. Klinik notların temsili için kullanılan ClinicalBert modeline ise buradan ⁸ eris,ilebilmektedir. Yapılan deneyler, 11GB kapasite bellekli NVIDIA RTX 2080 Ti ekran kartlıbilgisayarda kos,turulmuş,tır.

6.3 Deneysel Sonuçlar

Klinik notların, mortalite ve YBÜ'de kalma süresi tahmin problemlerine etkisini aras,tırmak için yapılan deney sonuçlarıÇizelge 6.4'de okuyucu ile paylas,ılmış,tır. Sonuçlardan görüleceği üzere, hastaya ait zaman serisi özniteliklerine ek olarak klinik notların problem çözümünde kullanılmasının AUROC, AUPRC ve F1 metrikler cinsinden pozitif etkisi olmuştur. Zamana bağlıöznitelikler GRU yöntemi ile temsil edilirken, klinik notlar SBERT yöntemi ile temsil edilmiş, ve birles,tirilen temsiller 2 katmanlıbir yapay sinir ağına girdi olarak verilmiştir.

Sonuçlar incelendiğinde mortalite problemi için AUROC ve AUPRC skorlarının %1 ile %2 arasında iyileştiği gözlemlenmektedir. F1 skorunda ise hastanede mortalite probleminde %4 başarımlarınartışığıgörülmüş,ken, yoğun bakımda mortalite probleminde %1.5 civarında olmuştur. Yoğun bakımda kalma süresi problemlerinde de AUROC ve AUPRC

⁸https://huggingface.co/emilyalsentzer/Bio_ClinicalBERT



Şekil 6.1: Mortalite ve Yoğun Bakımda Kalma Süresi temelli dört ayırklinik problemin tahmini için önerilen model mimarisi.

metrikleri bazında yaklaşık %1.5 iyileşme gözlemlenmektedir. Önerilen yöntemler içerisinde başarımın artış göstermediği tek metrik yoğun bakımda kalma süresinin 7 günden büyük olmasını tahmin etme problemindeki F1 skoru olmuştur. Bu durumun sebebinin, YBÜ'de 7 günden fazla kalmayı tahmin etme klinik görevinin sahip olduğu, yüksek veri dengesizliği probleminden dolayı olduğu düşünülmektedir (%7.9 oranında 7 günden fazla yoğun bakımda kalan hasta mevcuttur).

Çizelge 6.3: Çeşitli medikal kısaltmalar ve bu kısaltmaların dönüştürülen halleri.

Kısaltma	Dönüştürülen Hali
p.o.	orally
q.d.	once a day
i.m.	intramuscular
5x	a day five times a day
t.i.d	three times a day

Çizelge 6.4: Önerilen model performansının, temel olarak alınan modellerle karşılaştırılması. Her problem ve metrik için en yüksek skorlar vurgulanmıştır.

Görev	Yöntem	AUROC	AUPRC	F1
Hastanede Mortalite	GRU	87.36±0.003	51.25±0.007	41.28±0.033
	Önerilen Model	88.43±0.003	53.10±0.006	45.06±0.02
Yoğun Bakımda Mortalite	GRU	88.35±0.003	48.12±0.021	41.30±0.025
	Önerilen Model	89.00±0.003	49.68±0.008	42.54±0.027
Yoğun Bakımda Kalma>3 Gün	GRU	69.63±0.002	63.68±0.003	54.4±0.01
	Önerilen Model	70.25±0.004	64.96±0.004	55.15±0.014
Yoğun Bakımda Kalma>7 Gün	GRU	73.54±0.004	19.74±0.004	2.53±0.014
	Önerilen Model	75.14±0.002	21.35±0.004	0.0±0.000

6.4 Değerlendirme

Yapılan deney sonuçlarından da görüleceği üzere hastaya ait klinik rapor ve klinik personelin hasta için yazdığı notlar, mortalite ve YBÜ'de kalma süresinin tahmini problemleri için önemli bir veri kaynağı olmaktadır. Yaşamsal gözlem verileri veya laboratuvar sonuçları haricinde doktorun düşüncelerini aktardığı bu notlar, kıymetli bilgiler içermekle beraber hastanın geleceği hakkında önemli ipuçları taşımaktadır. Bu sebeple, klinik notların mortalite ve YBÜ'de kalma süresi problemleri dahil olmak üzere diğer klinik problemlerde de kullanılması gereken önemli bir veri kaynağı olduğu düşünülmektedir. Bu çalışmada, klinik notların temsili için önceden eğitilmiş ClinicalBERT modeli kullanılmış, ve SBERT yöntemi ile cümle bazında temsiller elde edilmiştir. Bu temsillerin ortalaması alınarak, her hastanın klinik notları temsil edilmiştir. Önerilen derin öğrenme tabanlı çok-kipli mimari ile klinik notlar ve zaman serisi öznitelikler basarılı bir şekilde işlenmiş, ve üzerinde çalışılan klinik problemler için farklı metrik türlerinde daha iyi sonuçlar alınmıştır. Gelecek çalışmalarda, klinik notların

zaman damgasının kesin olarak tutulduđu ESK verilerinde klinik notlar arasındaki zaman ilişkisini de kullanan derin öğrenme mimarileri önerilerek, model başarımlarının daha da arttırılmaya çalışılması planlanmaktadır.





7. ZAMAN SERİSİ ve İLAÇ TEMSİLLERİ ile TAHMİN ETME

7.1 Motivasyon

Derin öğrenme tabanlı modeller ile klinik problemlere çözümler üretmeye çalışan çalışmalar genellikle hastaya ait yaş, amsal gözlem verileri olan zaman serisi öznitelikleri kullanılmaktadır. Buna ek olarak, literatürdeki diğer çalışmalarda hastaya ait klinik notları veya klinik not içerisinden çıkartılan medikal terimleri yapay öğrenme tabanlı yöntemlere girdi olarak vererek çok-kipli yöntemler önermektedirler [70, 71]. Hastalara ait bir diğer önemli veri türü ise, hastaların YBÜ'de kaldıkları süre boyunca hastaya verilen ilaç bilgileridir. İlaç bilgisi ve yapısı özellikle kemi informatik (Cheminformatics), hesaplamalı biyoloji (computational biology), eczacılık (pharmaceutical) gibi alanlarda oldukça sık kullanılmasına rağmen ESK verileri ile üretilen çözümlerde şu ana kadarki bilginiz ışığında kullanılmamıştır. Bu çalışmamızda hastaya ait zaman serisi verileri, ilaç bilgilerinin moleküler yapısıyla birleştirilerek kullanılmıştır. Diğer çalışmalarımızda olduğu gibi bu çalışmada da YBÜ'de yatan hastaların ilk 24 saatteki verileri kullanılmış, ve hastanın YBÜ ve hastanede mortalite olma ihtimalleri tahmin edilmiş, ayrıca YBÜ'de kalma süresinin 3 ve 7 günden fazla olup olmayacağı problemleri çözülmeye çalışılmıştır.

Hastaya ait zaman-serisi öznitelikleri, Bölüm 3.2'de detayları anlatılan ve daha önceki çalışmalarda da kullanılan 104 öznitelikten oluşmaktadır. Hastaya ait ilaçlar ise MIMIC-III veri seti içerisindeki reçete (Prescription.csv) tablosundan çıkartılmıştır. İlaç isimlerinin yazımlarındaki farklılıkların düzeltilmesi için ve düzeltilen bu ilaç isimleri üzerinden bu ilaçların Simplified Molecular Input Line Entry System (SMILES) temsil karsılıklarının bulunabilmesi için ise çalışma kapsamında bir yöntem geliştirilmiştir. Hastaya ait ilaçlar, SMILES formatına dönüşürüldükten sonra vektörel forma dönüşürülmesi için Extended-Connectivity Fingerprints (ECFP) [58], Molecular Access System (MACCS) [59], Mol2Vec [60], ve Smiles-Transformer [61] yöntemleri olmak üzere 4 farklı şekilde temsil edilmiştir. Bu yöntemlerin detayları Bölüm 4.4'de okuyucu ile paylaşılmıştır. Son olarak ise, derin öğrenme modellerindeki klasik problemlerden biri olan açıklanabilirlik konusuna katkı sağlayabilmek için Bölüm 4.5'de detayları anlatılan SHapley Additive exPlanations (SHAP) yöntemi uygulanmıştır. Açıklanabilirlik ifadesi,

önerilen modelin tahminlerinde hangi özniteliklerden ne kadar yararlandığı bilgisinin elde edilmesini ifade etmektedir. Uygulanan bu yöntem sayesinde, modelin yapmış olduğu tahminler ile kullanılan öznitelikler arasındaki ilişki daha derin bir şekilde incelenebilmiş, ve hangi özniteliklerin modelin tahminine daha çok etkisi olduğu gözlemlenebilmiştir. Yöntemin çıktısı analiz edilerek hastanede mortalite problemi için önemli zaman-serisi öznitelikleri ve önemli ilaçlar tespit edilmeye çalışılmıştır.

7.2 Önerilen Yöntem

Tez kapsamında yapılan diğer deneylerde olduğu gibi bu çalışmada da MIMIC-III veri seti kullanılmıştır. Çizelge 7.1’de kullanılan veri setinin özet istatistikleri paylaşılmıştır. Dört farklı sınıflandırma problemi (hastane içi mortalite, YBÜ’de mortalite, YBÜ’de 3 günden fazla kalma, ve YBÜ’de 7 günden fazla kalma) tahmin edilmeye çalışılmıştır. Hastalara ait 104 adet zaman-serisi (yaş, amsal gözlem verileri ve laboratuvar sonuçlarını içeren) öznitelik Bölüm 3.2’de anlatıldı üzere çıkartılmıştır. İlaç isimlerinin bulunması ve dönüşümü. Çalışma kapsamında önerilen modellerde, zaman-serisi özniteliklerine (yaş, amsal gözlem verileri ve laboratuvar sonuçları) ek olarak klinik ilaç bilgileri de kullanılmıştır. MIMIC-Extract çalışması, ham MIMIC-III veri seti içerisinden zaman-serisi özniteliklerini çıkartmasına rağmen ilaç bilgilerini hastalar ile eşleştirmemektedir. MIMIC-III veri seti içerisinde, "Prescription" tablosu, hastalara tedavi amaçlı verilen ilaçların, tipi, ismi, gücü, dozu ve benzeri bilgileri saklamaktadır. Bu tablo içerisindeki verilerin örnek ve küçük bir kısmı Çizelge 7.2’de paylaşılmıştır.

İlaç isimlerini Prescription tablosu içerisinden çıkartabilmek ve elde edilen ilaç isimlerini vektörel hale dönüştürebilmek adına Şekil 7.1’de paylaşılan yöntem önerilmiştir. Bu yöntemde, ilk olarak hastalar ile hastaların ilk 24 saatte kullandıkları ilaç bilgileri eşleştirilmiştir. MIMIC-III içerisinde 2,255 tekil ilaç ismi ve 4,156,450 tane reçete kaydı olmasına rağmen, çalışma kapsamında kullanılan derlem içerisinde 1,967 tekil ilaç ve 592,946 adet reçete kaydı olmuştur. Elde edilen bu ilaç isimlerinin kirli/gürültülü olması sebebiyle çeşitli düzenli ifadeler (regular expressions), kural tabanlı isimler ve veri ön işleme adımları uygulanarak ilaç isimleri temizlenmeye ve tekilleştirilmeye çalışılmıştır. Örnek gürültülü ilaç isimleri ve bu ilaç isimlerine veri ön işleme adımları uygulandıktan sonraki durumları Çizelge 7.3’de paylaşılmıştır.

İlaç isimleri, veri ön işlemesinden sonra, PubChem [131] internet sitesi içerisinde aranarak, ilgili ilaçların SMILES kısımlarını elde edilmeye çalışılmıştır. PubChem ilk olarak 2004 yılında Amerika Birleşik Devletleri Ulusal Sağlık Enstitüleri (National Institutes of Health, NIH) içerisindeki Moleküler Kütüphaneler Programı (Molecular Libraries Program, MLP) tarafından yayınlanmıştır. Ağustos 2018 itibarıyla, PubChem,

Çizelge 7.1: MIMIC-III ve bu çalışmada kullanılan veri setinin istatistikleri.

	Hasta Sayısı	Hastane Bas_vuru Sayısı	YBÜ Bas_vuru Sayısı
MIMIC-III	46,520	58,976	61,532
MIMIC-III (>15 yaş,ından büyükler)	38,597	49,785	53,423
MIMIC-Extract	34,472	34,472	34,472
MIMIC-Extract (en az 24 saat YBÜ'de kalan hastalar)	23,937	23,937	23,937
Kullanılan Veri Seti(Medikal terimi olmayan hastaların elenmesinden sonra)	21,690	21,690	21,690

40 ülkeden 629 veri kaynağının katkıda bulunduğu, içerisinde 247.3 milyon madde tanımı, 96.5 milyon benzersiz kimyasal yapı, 10.000'den fazla hedef protein dizisi, 1.25 milyon biyolojik tahlil, ve 237 milyon biyoaktivite test sonucu içeren büyük veri tabanı olmuştur. Hastalara ait ilaç isimleri, açık kaynak olan Pubchempy⁹ kütüphanesi yardımı ile PubChem veri tabanında aranarak ilaca ait kimyasal yapı bilgisi elde edilmeye çalışılmıştır. Pubchempy kütüphanesi, PubChem veri tabanını servislerini çağırma işlevi sağlayan bir Python kütüphanesidir. MIMIC-III içerisindeki "Prescription.csv" içerisinde çıkarılan ve önerilen yöntemden geçirilerek temizlenen ilaç isimleri, Pubchempy kütüphanesi aracılığıyla aranarak PubChem veri tabanında aranmıştır. PubChem veri tabanı, aranan ilaçların moleküler yapısı, formülü, benzerleri, molekül ağırlıkları vb. daha birçok detaylı bilgi kullanıcılarına sunmaktadır.

Bu çalışmada, aranan ilaç ismine karşılık ilaçların yapısını satır notasyonu (line notation) ile temsil eden SMILES bilgisi eslesştirilmiştir. İlaç bilgileri, modeller içerisinde SMILES notasyonlarının ham hali ile kullanılmak yerine, Bölüm 4.4'de detayları anlatılan ECFP, MACCS, Mol2Vec ve Smiles-Transformer dahil olmak üzere dört farklı klinik ilaç temsiline dönüştürülerek kullanılmıştır. Bir ilacın SMILES verileri üzerinden ECFP ve MACCS ve Mol2Vec temsillerinin elde edilebilmesi için DeepChem [82]

⁹<https://pubchempy.readthedocs.io/en/latest/>

Çizelge 7.2: Reçeteli ilaçlar (Prescription) tablosundaki örnek ilaç ismi, genel ismi ve NDC (National Drug Code) örneği.

İlaç İsmi	Genel İlaç İsmi	NDC
Heparin	Heparin Sodium	63323026201.0
Acetaminophen	Acetaminophen	182844789.0
Lorazepam	Lorazepam	594091985307.0
Morphine Sulfate	Morphine Sulfate (Syringe)	409176230.0

Çizelge 7.3: Örnek klinik ilaç isimleri ve veri ön işleme adımlarından sonraki versiyonu.

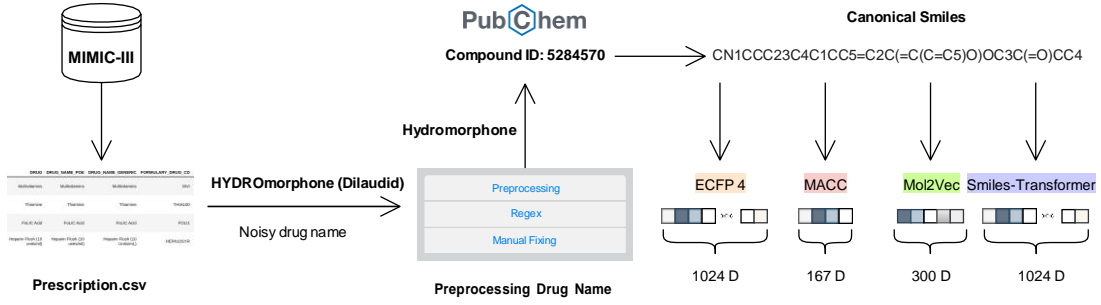
Gürültülü İlaç İsmi	Veri ön işlemesinden sonra
Heparin Flush (1000 units/mL)	Heparin Flush
NF* Rocuronium	Rocuronium
Methylene Blue 1%	Methylene Blue
23.4% Sodium Chloride	Sodium Chloride

kütüphanesi kullanılmıştır. Smiles-Transformer temsillerinin dönüşümü için ise ilgili çalışmanın [61] kendi Github hesabındaki¹⁰ kodlar kullanılmıştır. Deneylerde farklı molekül tanımlama yöntemlerinin kullanılmasının iki ana sebebi bulunmaktadır. İlk olarak, her moleküler tanımlayıcı, moleküllerin farklı yönlerini yakalayabilmektedir. İkinci olarak ise farklı moleküler temsil yöntemleri kullanılarak, hem bu temsil yöntemleri arasında kıyaslama yapmak, hem de ilaç molekül temsili kullanmanın başarımlara etkisinden emin olunmak istenmektedir.

Bu adımlardan sonra, hastaya ait hem zaman-serisi öznitelikleri hem de ilaç bilgilerinin vektörel formu elde edilmiştir. Model eğitimleri esnasında veri %70 eğitim, %10 validasyon ve %20 test verisi olacak şekilde bölünerek eğitimler gerçekleştirilmiştir. Çizelge 7.1’de çalışmada kullanılan veri sayısı ve Çizelge 7.4’de de her probleme ait örnek sayısı ve etiket dağılımı açıklanmıştır.

Zaman-serisi Öznitelikleri ile Eğitim. Hastaya ait saatlik olarak gruplanmış 104 adet zaman-serisi özneliğini modellere girdi olarak verebilmek adına Bölüm 4.1’de detaylarından bahsedilen Uzun kısa süreli bellek (Long short-term memory, LSTM) ve Geçitli tekrarlayan birim (Gated Recurrent Unit-GRU) yöntemleri kullanılmıştır. LSTM ağları, Tekrarlayan sinir ağları (Recurrent Neural Network, RNN) yönteminin bir varyantı olarak ortaya çıkmıştır. LSTM mimarisi içerisinde, RNN mimarisine ek olarak girdi (input), unutma (forget), ve çıkış, (output) kapıları bulunmaktadır. Bu sayede LSTM mimarileri önemli bilgileri çok daha uzun süre taşıyabilmekte ve böylece uzun vadeli ilişkileri yakalayabilmektedir. İkinci olarak ise LSTM yöntemi, RNN yönteminin problem yaşadığı kaybolan gradyan problemini (vanishing gradient problem), içerisindeki giriş ve unutma kapıları sayesinde yaşamamaktadır. RNN modelinin bir diğer varyantı olan ve LSTM modeline göre daha sade bir mimarisi olan GRU yönteminde ise reset ve güncelleme (update) kapıları bulunmaktadır. GRU modelinin LSTM modeline göre daha basit mimaride olması sebebiyle hesaplama maliyeti, LSTM mimarisine göre

¹⁰<https://github.com/DSPsleeporg/smiles-transformer>



Şekil 7.1: MIMIC-III içerisinde klinik ilaç isimlerini çıkartan yöntem. Çıkarılan ve ön işlemden geçirilen klinik ilaçlar PubChem içerisinde bulunduktan sonra farklı moleküler temsillere dönüşürmektedir.

daha verimlidir. Bu sebeplerden ötürü zaman-serisi öznitelikleri ile yapılan deneylerde LSTM ve GRU yapıları kullanılmış, olup, GRU modelinin LSTM mimarisine göre hem az da olsa daha başarılı sonuçlar vermesi hem de daha hızlı eğitim süresine sahip olmasından ötürü GRU yöntemi zaman-serisi özniteliklerini işlemek için seçilmiştir. Önerilen mimari de 128 nöron içeren tek katmanlı bir GRU kullanılmıştır.

Önerilen Çok-kipli Yöntem. Model mimarisine girdi olarak verilecek veri türlerinin farklılığından dolayı (zaman-serisi öznitelikleri ve ilaç temsilleri), Şekil 7.2’de gösterilen çok-kipli derin bir mimari önerilmiştir. GRU tabanlı ağlar, zaman-serisi öznitelikleri girdi olarak alıp 128-boyutlu bir öznitelik vektörü üretmektedirler. Hastaya ait ilaç vektörleri üzerinden öznitelik vektörü çıkartabilmek için ise Bölüm 4.1.4’de detayları anlatılan 1D Evrişimsel Sinir Ağı (1D CNN) mimarisi kullanılmıştır. D sembolü hastaya ait ilaç sayısı olarak kabul edilirse, ilaçların vektörel temsilleri $d = (d_1, d_2, \dots, d_D)$ olarak temsil edilebilmektedir. Hastaya ait ilaç sayısı da değişken olabileceğinden, çeşitli deneyler sonucundan her hastaya ait ilaç sayısı 64 olarak belirlenmiştir. Hastaya ait 64 ilaç bilgisi olmadığı durumlarda, ilaç temsil vektör boyutu kadar sıfırlardan oluşan bir vektör ile dolgu (padding) işlemi uygulanmıştır. $D \times 64$ boyutunda bir hasta ilaç temsili matrisi üzerine 1D CNN yöntemi uygulanarak 128 boyutlu bir öznitelik vektörü yaratılmaktadır. Evrişimsel sinir ağlarının ilk katmanları genellikle düşük-seviye öznitelikleri tespit etmek ve ileriki seviye katmanları ise daha detaylı öznitelikleri ortaya çıkartmak için tasarlanmasından ötürü [132], model mimarisi içerisinde, ard arda yerleştirilmiş, üç tane bir boyutlu evrişimsel katman eklemiş, ve filtre sayıları 32, 64, ve 128 olarak seçilmiştir. Çekirdek (kernel) büyüklüğü sabit olarak tutulmuş ve bütün evrişimsel katmanlar için 3 olarak seçilmiştir. Evrişimsel katmanların sonucundan sonra maksimum-örnekleme (max-pooling) işlemi gerçekleştirilerek 128-boyutlu öznitelik vektörü elde edilmiştir. Hasta-ilaç temsili matrisi üzerinden elde edilen 128-boyutlu öznitelik vektörü ile zaman-serisi öznitelikleri üzerinden elde edilen 128-boyutlu öznitelik vektörü birleştirilerek 3 katmanlı bir yapay sinir ağına girdi olarak verilmiştir. Yapay sinir ağı katmanlarının

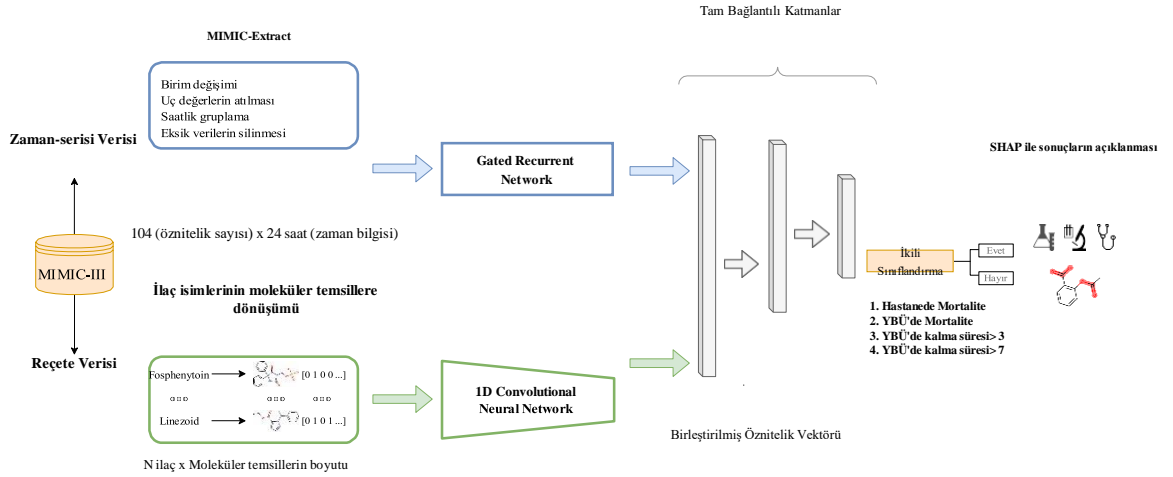
Çizelge 7.4: Klinik problemler için etiket dağılımları.

Problem İsmi	Örnek Sayısı	Etiket Dağılımı
Hastanede Mortalite	21,690	%10.37
YBÜ'de Mortalite	21,690	%6.98
YBÜ'de Kalma Süresi>3	21,690	%42.95
YBÜ'de Kalma Süresi>7 Gün	21,690	%7.69

nöron sayısı sırasıyla 1024, 512 ve 256 olarak belirlenmiştir.

Önerilen modelin son aşamasında bulunan tam bağlantılı katmanlarda (fully connected layers), aşırı öğrenme problemi (overfitting) engelleyebilmek için 0.3 oranında dropout [126] kullanılmıştır. Son katman haricindeki katmanlarda ise Rectified Linear Unit (ReLU) [127] aktivasyon fonksiyonu kullanılmıştır. ReLU aktivasyon fonksiyonu, $g(z) = \max\{0, z\}$ denklemine basitçe görüldüğü üzere, girdiyi 0 ile $+\infty$ arasında bir değere eşleştirmektedir. Tam bağlantılı katmanların son katmanında ise ikili sınıflandırma işlemini gerçekleştirebilmesi için sigmoid aktivasyon fonksiyonu kullanılmıştır. Önerilen model mimarisi, parametreler ve hiper parametreler, gerekli optimizasyon işleminden sonra belirlenmiş ve son haline getirilmiştir. Önerilen modelde kullanılan geri kalan hiper parametreler ise şu şekildedir: Yığın büyüklüğü (batch size) 32, regülasyon yöntemi olarak L2 ve oran olarak 0.05 seçilmiştir. Optimizasyon yöntemi olarak ise, 10^{-3} öğrenme oranı ve 10^{-2} sönüm (decay) oranı ile Adam [128] seçilmiştir. Bütün deneyler 100 iterasyon (epoch) olacak şekilde çalıştırılmıştır. Aşırı öğrenmeyi engellemek için ise erken durdurma (early stopping) yöntemi validasyon hatası (validation loss) üzerinde gözlemlenmiştir. Her problem için aynı modeller 10 farklı başlatma değeri (initialization seed) ile denenmiş ve ortalama değerler okuyucu ile paylaşılmıştır.

SHapley Additive exPlanations (SHAP) Yöntemi. Yapay ve derin öğrenme modellerinin açıklanabilirlik seviyesinin düşük olması ve kapalı-kutu yapıları sebebiyle, önerilen modellerin yaptıkları tahminin hangi sebebe dayanarak yaptığını anlamlandırmak zorlu bir konudur. Yapay zeka modellerinin açıklanabilir olması ise sistemi kullanan kişilerin, model tahminlerine olan güveni kazanabilmesi için önemli bir konudur. Geçmiş yıllarda araştırmacılar, önerdikleri modellerin açıklanabilir olması adına, yapısı gereği açıklanabilir olan karar ağaçları (decision-tree), lineer regresyon, lojistik regresyon gibi modeller kullanmışlardır. Güncel gelişmelerle beraber ortaya çıkan açıklanabilir yapay zeka (explainable artificial intelligence, XAI) alanında, herhangi bir modelin açıklanabilirliğini model-bağımsız (model agnostic) olarak yapabilen ve detayları Bölüm 4.5'de anlatılan, Local Interpretable Model-Agnostic Explanations (LIME) [133] veya SHAP [75] gibi

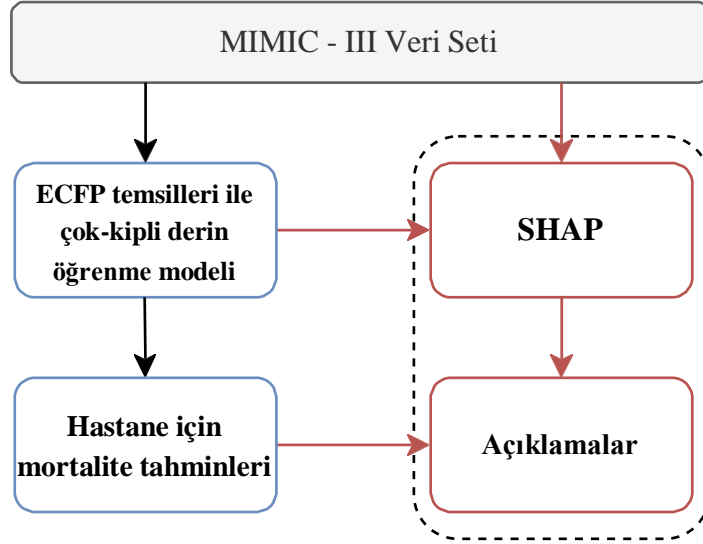


Şekil 7.2: Önerilen çok-kipli model mimarisinin özeti. İlk olarak, zaman-serisi ve ilaç verileri ön işlemden geçirilmiştir. İkinci olarak, öznitelik çıkarımı için, zaman-serisi özniteliklere GRU modeli, klinik ilaçlara ise 1D CNN modeli uygulanmıştır. Daha sonra ise, çıkartılan öznitelikler birleştirilerek yapay sinir ağına verilerek, 4 farklı klinik problem için tahmin yapılmıştır. Son olarak ise, SHAP yöntemi kullanılarak, hastane içi mortalite problemi ile önemli zaman-serisi öznitelikler ve ilaçlar arasında ilişki yakalanmaya çalışılmıştır.

yöntemler önerilmiştir. SHAP, 2017 yılında Lundber ve Lee tarafından her bir öznitelik için model tahminlerine etkisini ölçmek için ortaya çıkarılmış bir yöntemdir. Şekil 7.3'de gösterildiği üzere bu çalışmada, önerilen çok-kipli modele, SHAP yöntemi uygulanarak hastanede mortalite olma problemi için önemli zaman-serisi öznitelikleri ve klinik ilaç isimleri tespit edilmeye çalışılmıştır. SHAP yönteminin çıktıları, klinik uzmanlara ek bilgiler sağlayarak, onların ilgilendikleri hastaların hangi yaşamsal gözlem verilerine, laboratuvar sonuçlarına veya klinik ilaçlarına dikkat etmeleri gerektiğine yardımcı olabilecektir. Bildiğimiz kadarıyla gerçekleştirmiş olduğumuz bu çalışma, ESK alanında hem zaman-serisi özniteliklere hem de ilaçların moleküler temsillerine SHAP yöntemini uygulayan literatürdeki ilk çalışmadır.

7.3 Deneysel Sonuçlar

Deney Ayarları. Üzerinde çalışılan bütün klinik problemler MIMIC-Extract [87] çalışmasında tartışıldığı üzere birbirinden ayırda dört ikili sınıflandırma problemi olarak ele alınmıştır. Bu dört problemten üç tanesinde (LOS > 3 hariç diğer problemlerde) sınıf dengesizliği problemi yaşanmaktadır. Bu durum, bu problemi yaşayan klinik görevlerde, azınlık olarak bulunan sınıfın tahmin edilebilirliğini zorlaştırmaktadır. Bu problem, ilgili klinik görevlerde, sınıf ağırlıkları, sınıfların bulunma sıklığı ile ters orantılı olacak şekilde ayarlanarak çözülmüştür. Modellerin eğitilmesi, parametre optimizasyonu ve test edilmesi için veri seti %70, %10, ve %20 oranlarında eğitim, validasyon ve test



Şekil 7.3: SHAP uygulamasının çalışma içerisindeki kullanımı.

seti olarak ayrılmıştır. Model eğitimlerinin değerlendirilmesi aşamasında detayları Bölüm 4.6’de anlatılan, F1 skoru, AUROC ve AUPRC metrikleri kullanılmıştır.

Model eğitimlerinde kullanılan derin öğrenme tabanlı yöntemler arka planda Tensorflow [122] olacak şekilde Keras [123] kütüphanesi ile geliştirilmiştir. İlaçların isimlerinin çıkartılması ve temsillerinin elde edilmesi için ise Pubchempy¹¹, RDkit¹², ve Deepchem [82] çalışmalarından yararlanılmıştır. Bütün modeller Intel i7-9700K CPU ve NVIDIA RTX 2080 Ti ekran kartı mevcut olan donanım üzerinde eğitilip test edilmiştir. İlgili çalışma kodları ve adreslerinden^{13,14} erişilebilmektedir.

Önerilen Yöntemin Deney Sonuçları. Bu bölümde yapılan deney sonuçları değerlendirilmiştir. Bütün klinik problemler aynı eğitim verisi ile eğitilip, aynı test verisi ile test edilmiştir. Sadece zaman-serisi ile eğitilen modellerde, MIMIC-Extract çalışmasının çıktısı olarak elde edilen 104 adet zaman-serisi öznitelik kullanılmıştır. Bu öznitelikler ile LSTM ve GRU modelleri eğitilmiştir. Sonuçlar incelendiğinde, GRU daha basit bir mimariye sahip olmakla beraber LSTM modellerine göre az da olsa (ortalama %0.5 - %1) daha iyi sonuç göstermiştir. Çizelge 7.5’de görüldüğü üzere, GRU yöntemi hastane içi mortalite tahmininde 84.73, YBÜ’de mortalite tahmininde 86.91, YBÜ’de Kalma Süresi > 3 probleminde 69.09 ve YBÜ’de Kalma Süresi > 7 probleminde 67.52 AUROC performansı göstermiştir.

Bu sonuçları iyileştirmek adına, klinik ilaçların temsilleri modele dahil edilerek yeni deneysel çalışmalar gerçekleştirilmiştir. İlaç temsilleri üzerinden öznitelik çıkarımı

¹¹<https://pubchempy.readthedocs.io/>

¹²<https://www.rdkit.org>

¹³<https://github.com/tanlab/MIMIC-III-Clinical-Drug-Representations>

¹⁴https://github.com/tanlab/Explainable_MIMIC_III

yapmak için 1D CNN yöntemi kullanılmıştır. Sonuçlar incelendiğinde, ilaç bilgisinin kullanılmasının üzerinde çalışılan bütün problemlere olumlu etki sağladığı gözlemlenmektedir. Dört farklı ilaç temsilinin de model başarımlarına pozitif etkisi olmasına karşın, en fazla iyileşmenin ECFP temsil yöntemi ile olduğu görülmektedir. Çizelge 7.5'de incelendiğinde, hastanede mortalite probleminde performanslar AUROC metriği cinsinden %3, AUPRC olarak %8, ve F1 skoru olarak ise %4 olarak iyileşmiştir. YBÜ'de mortalite probleminde, hastanede mortalite göre iyileşme biraz daha düşük olmasına rağmen sırasıyla AUROC, AUPRC ve F1 metriklerinde %2, %5, %3 oranında gerçekleşmiştir.

YBÜ'de Kalma Süresi > 3 problemi sonuçları incelendiğinde, 72.23 AUROC skoru ile ECFP diğer bütün yöntemlerden daha başarılı olmuştur. Ayrıca sadece zaman-serisi öznitelikleri ile yapılan deney sonucuna göre AUROC metriği bazında %3'lük bir iyileştirme gözlemlenmektedir. AUPRC ve F1 skoru incelendiğinde de ECFP skorlarının, zaman-serisi başarımlarına göre yaklaşık olarak %3'lük bir iyileşme gösterdiği gözlemlenmektedir. YBÜ'de Kalma Süresi > 7 problemi incelendiğinde ise ilaç bilgisinin kullanılmasının model başarımına çok daha büyük bir etkisi olduğu gözlemlenmiştir. AUROC metriği türünden ECFP temsil yöntemi ile yapılan deneylerde %74.71 başarımları yakalanmıştır. Bu oran, zaman-serisi deney sonucundan yaklaşık olarak %7 daha iyi bir sonuç olmuştur. İlaç temsil yöntemleri kendi aralarında karşılaştırıldığında ECFP en iyi sonucu vermektedir. İkinci olarak MACCS yöntemi en iyi sonucu vermesine rağmen, ortalama olarak ECFP yöntemi MACCS, Mol2Vec, ve Smiles-Transformer yöntemlerine %0.5 ile %2 arasında fark attığı gözlemlenmektedir.

7.4 Değerlendirme

Bu bölümde, yapılan deneyler sonucunda elde edilen sonuçlar tartışılmıştır. Ayrıca hastanede mortalite olma problemi tahminlenirken kullanılan özniteliklerden önemli ve önemsiz olanları keşfetmek için modelin daha açıklanabilir olmasına çalışılmaya çalışılmıştır. Bu sebeple, modele açıklanabilirlik katmak adına SHAP yöntemi uygulanmıştır. Bu işlemler sonucunda hastanede mortalite problemi için eğitilen modelin, yaptığı tahminler esnasında hangi özniteliklerden çok yararlanıp hangilerinden az yararlandığı tespit edilmiştir. Elde edilen bu önemli ve önemsiz özniteliklerin gerçekten modele olan etkisini araştırmak için bu öznitelikler çıkartılarak çeşitli kez deneyler tekrarlanmıştır. Ayrıca hastanede mortalite problemi için bulunan önemli/önemsiz zaman-serisi öznitelikleri ve klinik ilaç isimleri okuyucu ile paylaşılmıştır.

Klinik İlaç Temsilleri Kullanmanın Model Performansına Etkisi. Üzerinde çalışılan dört klinik görevin (hastanede mortalite, YBÜ'de mortalite, YBÜ'de Kalma Süresi > 3, YBÜ'de Kalma Süresi > 7) tahmin performansını iyileştirmek için çok-kipli

Çizelge 7.5: Dört klinik problem için alınan basarımların ortalama sonuçları ve birbirleri ile karşılaştırılması.

Problem	Yöntem	İlaç Temsili	AUROC	AUPRC	F1
Hastanede Mortalite	GRU	-	84.73±0.010	46.24±0.027	42.29±0.015
		MACCS	87.43±0.003	52.01±0.009	45.09±0.008
	Önerilen Model	Mol2Vec	86.55±0.002	49.97±0.005	42.85±0.009
		Smiles-Transformer	86.67±0.002	49.57±0.004	43.15±0.007
		ECFP	87.80±0.003	53.24±0.007	46.55±0.015
YBÜ' de Mortalite	GRU	-	86.91±0.008	43.08±0.017	42.22±0.02
		MACCS	88.62±0.003	48.36±0.01	45.24±0.013
	Önerilen Model	Mol2Vec	87.59±0.003	45.21±0.008	42.67±0.006
		Smiles-Transformer	87.55±0.002	45.30±0.007	43.27±0.009
		ECFP	88.76±0.003	48.54±0.014	45.65±0.01
YBÜ' de Kalma Süresi > 3 Gün	GRU	-	69.09±0.004	62.90±0.004	57.44±0.009
		MACCS	71.27±0.004	64.56±0.003	58.90±0.009
	Önerilen Model	Mol2Vec	70.09±0.003	65.54±0.003	59.88±0.003
		Smiles-Transformer	69.49±0.002	62.68±0.002	58.35±0.006
		ECFP	72.23±0.003	65.69±0.003	60.88±0.003
YBÜ' de Kalma Süresi > 7 Gün	GRU	-	67.52±0.019	15.20±0.018	21.31±0.019
		MACCS	73.52±0.008	20.36±0.009	27.62±0.006
	Önerilen Model	Mol2Vec	73.29±0.005	20.09±0.002	25.69±0.004
		Smiles-Transformer	73.51±0.006	20.30±0.005	25.89±0.006
		ECFP	74.71±0.002	21.00±0.004	28.85±0.007

bir derin öğrenme mimarisi önerilmiştir. Zaman-serisi özniteliklerinin yanında ilaç temsil bilgilerinden verimli bir şekilde yararlanabilmek adına GRU, 1D CNN ve tam bağlantılı yapay sinir ağları beraber kullanılmıştır. Deneysel sonuçlardan görüleceği üzere klinik ilaçların zaman serisi öznitelikler ile beraber kullanıldığı deneyler en iyi sonuçları vermiştir. Hastanede mortalite problemi sonuçları incelendiğinde, farklı moleküler temsil yöntemleri ile eğitilen modellerin benzer AUROC sonucu gösterdiği görülmektedir. AUROC bazında iyileşmenin genel olarak ortalama %2.5 civarında olduğu görülürken, AUPRC ve F1 skoru bazında modeller arasında farklar olduğu görülmektedir. Mol2Vec ve Smiles-Transformer temsilleri, sadece zaman-serisi ile eğitilen modellere göre %4 AUPRC ve %0.3 F1 skoru daha iyi sonuç vermesine karşın ECFP ve MACCS temsillerinin gerisinde kalmışlardır. Yoğun bakımda yatan hastaların mortalite tahminlerinde, moleküler temsillerin model başarımlarına etkisi ve yarattıkları farklar hastanede mortalite tahmini modelindeki sonuçlara benzemektedir. Sonuçlar incelendiğinde, geleneksel molekül tanımlayıcı yöntemlerden MACCS ve ECFP temsilleri, yapay öğrenme tabanlı Mol2Vec ve Smiles-Transformer yönteminden daha iyi sonuç göstermektedirler. Sonuçlardan da görüleceği üzere, klinik ilaçların temsillerinin, zaman-serisi öznitelikler ile beraber kullanılmasında mortalite problemi başarımlarını arttırmaktadır. Aynı zamanda, moleküler temsillerin kendi aralarındaki kıyaslamada ise ECFP en başarılı sonuçları veren moleküler temsil olmakla beraber ECFP yöntemini MACCS izlemektedir.

Yoğun bakımda kalma süresinin 3 ve 7 günden fazla olup olamayacağı tahmin eden model sonuçları incelendiğinde, klinik ilaç bilgilerinin model sonuçlarına olumlu etki yarattığı görülmektedir. Üzerinde çalışılan dört klinik problem içerisinde sınıf dağılım problemini en az yaşayan YBÜ'de kalma süresi > 3 gün probleminde %60.89 F1 skoru alınarak, klinik problemler arasındaki en yüksek F1 skoru elde edilmiştir. Yoğun bakımda kalma süresi tahmini problemlerinde de en iyi sonuçlar ECFP temsili kullanıldığında alınmış ve ortalama olarak bütün metriklerde %3-%4'lük bir iyileştirme sağlamıştır. Son klinik problem olan YBÜ'de kalma süresi > 7 günde ise, çok yüksek bir başarımla iyileşmesi gözlemlenmiştir. İlaç temsillerinin kullanılması ve çok-kipli derin öğrenme modellerine geçişte AUROC ve AUPRC metriklerinde %6 olarak görülen iyileşme, F1 metriği için %7 olmuştur.

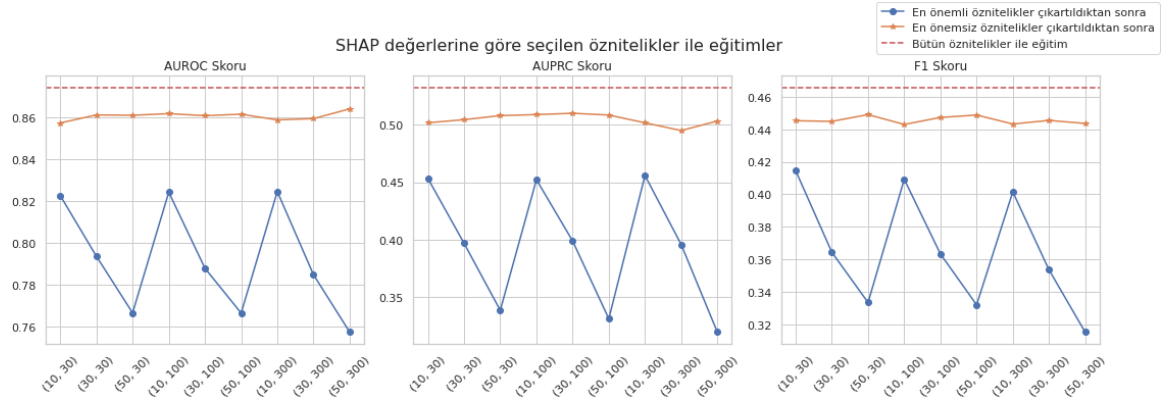
Özetle, deneysel sonuçlar incelendiğinde dört klinik problem içinde ilaçların moleküler temsillerinin kullanılmasının model başarımlarını arttırdığı görülmektedir. Yapılan deneyler sadece ilaç temsillerini kullanmanın avantajını göstermekle kalmamakta ayrıca farklı ilaç temsil yöntemlerinin de birbirleri ile kıyaslanmasına olanak sağlamaktadır. Alınan sonuçlar, ECFP ve MACCS yöntemlerinin açık bir şekilde önceden eğitilmiş olan Mol2Vec ve Smiles-Transformer yöntemlerinden daha iyi sonuç verdiğini göstermek-

tedir. Bunun ilk sebebinin Mol2Vec ve Smiles-Transformer yöntemlerinin MIMIC-III içerisindeki klinik ilaç verileri dışında başka molekül veri setleri ile eğitilmesi olarak düşünülmektedir. İkinci olarak ise, literatürde de oldukça sık bir şekilde kullanılan geleneksel moleküler tanımlayıcılar olan ECFP ve MACCS'un klinik ilaç yapılarını oldukça iyi temsil edebilmesi olarak düşünülmektedir.

SHAP sonuçlarına göre seçilen öznitelikler ile ek deneyler. Hastanede mortalite probleminin sonuçlarının açıklanabilirliğinin sağlanabilmesi için önerilen çok-kipli derin öğrenme tabanlı yöntem (ECFP temsillerinin kullanıldığı) Bölüm 4.5'de anlatılan SHAP yöntemi uygulanmıştır. Örnek sayısının fazlalığı ve model büyüklüğünden dolayı SHAP'ın bir varyasyonu olan GradientShap yöntemi kullanılmıştır. Bu yöntem, temel dağılımdan rastgele örnekleme yoluyla gradyan değerlerini hesaplayarak, orijinal SHAP değerlerine yaklaşan değerler üretmektedir. Zaman-serisi özniteliklerinin SHAP değerini hesaplamak için, özniteliklere ait 24 saatlik veri içerisindeki değerlerin mutlak değerleri toplanmıştır. Daha sonra ise tüm hastalara ait bu değerlerin ortalaması alınarak, en yüksek değere sahip zaman-serisi özneliğinden en düşüğe göre sıralama yapılmıştır. Klinik ilaçların önem sıralaması için ise, her ilaca ait 1024 boyutlu (ECFP temsili) vektördeki değerlerin mutlak toplamı alınmıştır. Daha sonra ise aynı ilaçların SHAP değerlerinin ortalaması alınarak her bir ilaç için önem puanı hesaplanmış ve sonrasında bu değere göre önem sıralaması gerçekleştirilmiştir.

Zaman-serisi özniteliklerin ve klinik ilaçların SHAP değerlerine göre sıralanmasından sonra, bu öznitelikler modelden çıkartılarak sistematik deneyler gerçekleştirilmiştir. Bu deneylerin amacı, SHAP yönteminin bulunduğu özniteliklerin gerçekten önemli olup olmadığını anlamaya çalışmak olmuştur. Bu deneyleri gerçekleştirebilmek için öncelikli olarak en önemli/önemsiz 10, 30, 50 zaman-serisi özneliği ve 30, 100 ve 300 klinik ilaç ismi seçilmiştir. Önemli özniteliklerle dokuz, önemsiz özniteliklerle de dokuz olmak üzere toplam 18 tane model eğitilmiştir. Şekil 7.4'de yapılan 18 deneyin sonuçları AUROC, AUPRC ve F1 skoru türünden paylaşılmıştır. Mavi noktalar ile temsil edilen deneyler, en önemli özniteliklerin çıkartılması ile yapılan deneylerin sonucunu gösterirken, turuncu noktalar ile temsil edilen deneyler ise en önemsiz özniteliklerin çıkartılması ile eğitilen modellerin başarımlarını göstermektedir. Kırmızı düz çizgi ise, hastanede mortalite probleminin çözümünde tüm özniteliklerin kullanıldığı deney sonucunu ifade etmektedir.

Şekil 7.4'den görüleceği üzere, SHAP önemli ve önemsiz öznitelikleri doğru bir şekilde ayırt edebilmiştir. Önemsiz öznitelikler çıkartılarak yapılan deneylerde de başarımlar belli miktarda düşmesine rağmen, önemli öznitelikler model içerisinden çıkartıldığında model başarımları düşme miktarı oldukça fazla olmaktadır. F1 skoruna göre, en önemli 10 zaman-serisi öznelik ile 30 klinik ilacın çıkartılarak yapıldığı deney sonucunda



Şekil 7.4: SHAP çıktılarıyla önemli özelliklerin seçimi. Deneysel sonuçlar üç farklı figürde gösterilmiştir. Figürlerdeki X-ekseni (x,y) formatında olmak üzere çıkartılan özellik sayılarını temsil etmektedir. x sembolü çıkartılan zaman-serisi özellik sayısını temsil ederken, y ise çıkartılan klinik ilaç sayısını temsil etmektedir.

Çizelge 7.6: Model çıktısına en çok ve en az katkı veren ilaç isimlerinin listesi.

Klinik İlaç İsimleri	
En Önemli İlaç İsimleri	En Önemsiz İlaç İsimleri
Meloxicam	Milk thistle
Guanine	Tipranavir
Vincristine	Trimethoprim
Doxycycline monohydrate	Precose
Oxaliplatin	Allopurinol sodium
Solifenacin succinate	Cromolyn
Fenoldopam mesylate	Delavirdine mesylate
Pristiq	N,6-dimethylhept
Sitagliptin	Didanosine
Zantac	Isosorbide

F1 skoru 46.5'den 41.4'e düşmektedir. Deney sonuçları incelendiğinde bu düşüşün ana sebebinin zaman-serisi özellikleri olduğu görülmektedir. Buna rağmen, klinik ilaçların da modelin başarımına katkıda bulunduğu gözlemlenebilmektedir. Örneğin, model eğitiminden çıkarılan zaman-serisi özneliğini sabit ve 10 olarak tutup, çıkarılan klinik ilaç sayısını (30, 100, 300) olarak değiştirdiğimizdeki deney sonuçları incelendiğinde, F1 skorunun 41.4, 40.8, 40.1 olarak azaldığı görülmektedir. Önemsiz olarak işaretlenen özelliklerin çıkartılmasıyla yapılan deneylerde ise böyle bir fark gözlemlenmemektedir.

Çizelge 7.7: Model çıktısına en çok ve en az katkııveren zaman-serisi öznitelik listesi.

Zaman Serisi Öznitelikleri	
En Önemli Öznitelikler	En Önemsiz Özniteliker
glasgow coma scale total	creatinine pleural
blood urea nitrogen	lymphocytes atypical csl
diastolic blood pressure	creatinine bodyfluid
anion gap	albumin pleural
tidal volume set	lymphocytes percent
mean corpuscular volume	albumin ascites
mean corpuscular hemoglobin concentration	creatinine ascites
sodium	calcium
creatinine	red blood cell count pleural
mean blood pressure	red blood cell count urine

Hastanede mortalite problemi için en önemli zaman-serisi öznitelikleri ve ilaç isimleri.SHAP yöntemi sayesinde klinik veri setleri içerisindeki önemli öznitelikleri anlamak için yeni bir yöntem geliştirilmiştir. Hastanede mortalite probleminin açıklanabilirliğini artırmak için uygulanan SHAP yöntemi ile, özniteliklerin önemi sıralanmış ve en önemli/önemsiz 10 zaman-serisi özniteliği ile klinik ilaç ismi Çizelge 7.6'de ve Çizelge 7.7'de okuyucu ile paylaşılmıştır. Mortalite ile ilişkili olabilecek önemli özniteliklerin bulunması, klinik uzmanların hastalara daha iyi bir tedavi sunabilmesine veya hastaları önceliklendirebilmelerine olanak sağlayabilecektir.

Zaman-serisi öznitelikleri içerisinde SHAP değeri en yüksek olanlar incelendiğinde, bu özniteliklerin daha çok yaşamsal gözlem verileri içerisindeki kritik hayati bulgular ile ilgili olduğu gözlemlenmektedir (Örneğin, glasgow koma skalası, diyastolik kan basıncı, anyon açığı). En yüksek SHAP değerine sahip zaman-serisi özniteliği olan "Glasgow Koma Skalası(GKS)" hastaların bilinç durumunu/bozukluğunu ölçmek ve değerlendirmek için kullanılmaktadır. GSK'ya göre ölçümlerde, hastaya 3 (derin bilinç kaybı) ile 15 arasında puan verilmektedir. Bu puanlama üç parametre dikkate alınarak gerçekleştirilir. Bu parametreler; gözler, sözlü yanıtlar ve motor tepkisinden oluşmaktadır. SHAP sonuçlarına göre bulunan en önemli ikinci öznitelik "Kan Üre Azotu (Blood Urea Nitrogen, BUN)" değeri olmuştur. Literatürde, birçok çalışma BUN değerini kullanarak mortalite üzerine çalışmalar gerçekleştirmiştir [134, 135]. "Diyastolik Kan Basıncı(Diastolic blood pressure, DBP)" kalbin atardamar duvarlarına uyguladığı kuvveti ifade etmektedir. Taylar vd. [136] yaptıkları çalışmada, özellikle 50 yaşından genç kişiler için DBP değerinin mortalite için önemli özniteliklerden biri olduğunu vurgulamaktadır. Zaman-serisi özniteliklerinin yanısıra klinik ilaçlarda SHAP

değerlerine göre sıralanmış ve Çizelge 7.7’da bu değerler paylaşılmıştır. İlaçların karmaşık yapısından ötürü, mortalite tahmini ile ilaçlar arasında ilişki kurmak ve yorumlamak zor bir problemdir. Ancak yine de bu bilgilerin doktorların tedavi süreçlerine katkı sağlayabileceğine, doktorların tedavi planlarını da geliştirerek ilaçların birbirilerini arasındaki etkileşimlerinden kaynaklanabilecek ciddi durumların önüne geçebileceği düşünülmektedir.





8. SONUÇ ve ÖNERİLER

Teknolojinin son yıllarda hızla gelişmesi ve her geçen gün sayısı artan dijital uygulamalar sayesinde, veriler çok çeşitli kaynaklardan toplanmakta ve veri miktarı katlanarak artış göstermektedir. Bu verileri doğru yöntemlerle işleyerek, katma değerli bilgi üretmek ise artık çok daha olasıdır. Sadece veri miktarının artışı ve teknolojinin gelişmesi değil, aynı zamanda yapay öğrenme, derin öğrenme veya veri bilimi gibi alanlarda çalışan araştırmacıların da artmasıyla beraber algoritmik gelişmelerde hızlanmıştır. Birçok farklı alan içerisindeki problemlerin çözümü için farklı farklı yapay zeka tabanlı yöntem ve algoritma önerilmiştir. Bu alanlar içerisinde sağlık alanı, insanlar var olduğundan beri, en önemli alanlardan biri olmuştur. Birçok alanda olduğu gibi sağlık alanında da dijitalleşmenin artmasıyla beraber sağlık alanında da veriler dijital ortama taşınmış ve yapay zeka modellerine girdi olabilecek noktaya gelmiştir. Sağlık alanında veri mahremiyeti problemi ise literatürde açık kaynak olarak paylaşılan ESK veri setleri ile aşılmıştır. Bu sayede yapay zeka ve sağlık alanında çalışan araştırmacılar, birçok klinik problemin çözümünü araştırmak için bu verilerden yararlanmakta ve yapay zeka yöntemleri ile çözümler aramaktadırlar.

Yapılan bu çalışmada, ESK veri seti içerisindeki birçok veri türü bir arada kullanılarak, yoğun bakımda yatan bir hastanın mortalite olup olmamasını ve yoğun bakımda kalma süresini tahmin edebilen modeller geliştirilmiştir. Geliştirilen bu modeller, literatürdeki popüler, açık kaynak, elektronik sağlık kayıt veri seti olan MIMIC-III kullanılarak eğitilmiştir. MIMIC-III veri setinin ham halinin doğrudan modellere girdi olarak verilememesinden ötürü, veri temizleme, birleştirme, öznitelik haline dönüştürme gibi işlemler gerçekleştirilmiştir. İşlenen bu veriler sonrasında çok-kipli derin öğrenme modellerine girdi olarak verilerek, mortalite ve yoğun bakımda kalma sürelerinin tahmini gerçekleştirilmiştir. Yapılan deneylerde, hastaya ait özniteliklerden hangilerinin model başarımına pozitif etki edeceği araştırılmış olup, aynı zamanda bu özniteliklerin nasıl bir derin öğrenme mimarisi ile daha efektif bir şekilde işlenebileceği konusunda çalışmalar gerçekleştirilmiştir.

İlk olarak, üzerinde çalışılacak olan problemler formüle edilmiştir. İki temel problem olan mortalite tahmini ve yoğun bakımda kalma süresini tahmin etme klinik problemleri dört farklı sınıflandırma problemi olarak tanımlanmıştır: Bu problemler; hastanede

mortalite olup olmama, yoğun bakımda mortalite olup olmama, yoğun bakımda 3 günden fazla kalıp kalmama ve yoğun bakımda 7 günden fazla kalıp kalmama tahmini olarak tasarlanmıştır. Bu problemleri çözmek için ise yoğun bakımda kalan hastaların ilk 24 saatlik verileri kullanılmıştır. Bu kapsamda gerçekleştirilen ilk çalışmada, hastaya ait 104 adet zamana bağlı öznitelik (yaşamsal gözlem verileri, laboratuvar sonuçları) ile beraber hasta için yazılan klinik notlar beraber kullanılmıştır. Yapılan bu deneylerde, klinik notlar, dönüş türücü tabanlı BERT modeli ile temsil edilmiştir. Deneylerin temel amacı, zaman-serisi özniteliklerin yanında klinik notların kullanmanın üzerinde çalışılan klinik problemlere etkisini gözlemlemek ve model başarımını iyileştirmektedir. Alınan sonuçlar incelendiğinde, klinik notların kullanmanın model başarımlarından mortalite ile ilgili olanlara F1 skoru bazında %4'lük bir iyileştirme, yoğun bakımda kalma süresi tahminine ise %1.5'lük bir iyileştirme sağladığı görülmüştür.

Ayarlanan ikinci deney düzeneğinde, yapılan ilk deneyler ile aynı klinik problemler üzerinde çalışmalar gerçekleştirilmiştir. Bu çalışmanın ilk çalışmadan farkı ise, zaman-serisi öznitelikler ile beraber kullanılan klinik notların doğrudan temsil edilmesi yerine, varlık isim tanıma yöntemi kullanılarak klinik notlar içerisinde çıkartılan medikal terimlerin kullanılması olmuştur. Çalışmada öncelikli olarak medikal terimlerin en başarılı nasıl temsil edilebileceği araştırılmıştır. Bunun için Word2Vec, FastText, Doc2Vec gibi birçok yöntem ile deneyler gerçekleştirilmiştir. Ardından bu temsiller üzerinden öznitelik vektörünün nasıl çıkartılması gerektiği ile ilgili araştırmalar gerçekleştirilerek ortalama alma, 1D CNN gibi yöntemler ile çok-kipli derin öğrenme tabanlı yöntemler denenmiştir. Alınan sonuçlar incelendiğinde, birçok deney için en iyi sonuç, medikal terimlerin Word2Vec ile temsil edilmesi ile alınmıştır. Öznitelik vektörü çıkartmak için ise, 1D CNN yöntemi en iyi sonucu vermiştir. Medikal terimleri, zaman serisi öznitelikler ile beraber kullanmanın mortalite problemlerindeki F1 skorunda yaklaşık olarak %5, yoğun bakımda kalma süresini tahmin etme probleminde ise %2'lik bir iyileştirme yakaladığı görülmektedir.

Son çalışmada da diğer iki çalışmadaki klinik problemlerin aynı üzerinde çalışılırken, farklı olarak, hastalara ait ilaç bilgilerinin moleküler yapısını kullanmanın, üzerinde çalışılan klinik problemleri tahmin etmesine etkisi araştırılmıştır. Hastaya verilen klinik ilaçların moleküler yapısının vektörel hale dönüştürülmesi için ECFP, MACCS, Mol2Vec, Smiles-Transformer yöntemlerinden yararlanılmıştır. Ek olarak, SHAP yöntemi kullanılarak modele açıklanabilirlik kazandırılmıştır. Bu sayede hastanede mortalite problemi için önemli olan yaşamsal gözlem verileri, laboratuvar sonuçları, klinik ilaç isimleri elde edilebilmiştir. Mortalite tahmini yapmak için önemli olduğu bulunan bu öznitelikler, modelden çıkartılarak tekrar modeller eğitilmiş ve bu özniteliklerin gerçekten modele olan etkisi ölçülerek, önemli bulunan özniteliklerin

önemi teyit edilmiştir. Aynı zamanda deneysel sonuçlar, hastaya ait ilaç bilgilerinin moleküler yapısının zaman-serisi öznitelikler ile beraber kullanılmasının klinik problem tahminlerine olumlu etki ettiğini göstermiştir. Mortalite problemleri için %3-4 oranında F1 skorunda iyileşme olurken, yoğun bakımda kalma süresinin 3 günden fazla olduğunu tahmin etme probleminde %3.5, 7 günden fazla olduğunu tahmin etmede %7 civarında gerçekleşmiştir.

Özetle, yapılan çalışmalarda literatürde sıklıkla kullanılan yasa,msal gözlem verileri ve laboratuvar sonuçlarına ek olarak, hastaya ait klinik notları, medikal terimleri ve ilaçların moleküler yapılarının kullanmanın önemi gösterilmiştir. Bu deneyler esnasında bu verilerin nasıl temsil edilmesi gerekliliği dört farklı klinik problem üzerinde denenmiştir. Ayrıca son çalışmada, modele açıklanabilirlik yeteneği eklenerek, hem bu modelleri geliştiren kişiler için bilgi verici çıktı elde edilmiş, hem de bu modelleri kullanabilecek klinik uzmanlara, modelin çıktısı olan mortalite ve yoğun bakımda kalma süresi üzerine tahminlerin sebebini ortaya çıkartmıştır.

8.1 Gelecek Çalışmalar için Öneriler

Sağlık alanının hızla gelişmesi, derin öğrenme yöntemlerinin çeşitliliği, elektronik sağlık kayıt verileri içerisindeki veri türlerinin zenginliği ve çok sayıda klinik problemin formülü edilebilmesinden ötürü gelecekte tez kapsamında yapılan çalışmalara ilave geliştirmeler yapılabilir durumdadır. Aşağıda çeşitli gelecek çalışma konularından bahsedilmiştir.

1. MIMIC-III veri seti haricinde bir başka veri seti kullanılarak, önerilen modeller tekrarlanabilir, parametreler kontrol edilebilir, ve önerilen bu yöntemler genelleştirilmeye çalışılabilir. Ayrıca önerilen bu çalışmanın, farklı ESK veri setleri arasında adaptasyon (domain adaptation) için yeni yöntemler geliştirilebilmenin önünü açabileceği düşünülmektedir.
2. Hastaların zaman-serisi özniteliklerinin, Gramiyen Açıs,al Alanlar (Gramian Angular Field, GAF) veya benzeri yöntemler ile görüntüye çevrilerek kullanılması ilginç yeni çalışmaların önünü açabilecektir.
3. Hastalara ait radyolojik görüntülerinin, hastanın mortalite olma ihtimalini tahmin etmede önemli bir yer tutabileceği düşünülmektedir. MIMIC-IV veri seti içerisindeki bu görüntüler, önerilen çok-kipli derin öğrenme modeline yeni bir veri türü olarak katılabilecek ve üzerinde çalışılan klinik problemlerin parametrelerini incelenebilecektir.

4. Tez kapsamında yapılan deneylerde kullanılan MIMIC-III veri setinin içerisindeki klinik notlara ve ilaç kullanım bilgilerine ait zaman damgaları yeterince hassas ve güvenilir olmamasından ötürü, bu veri türleri için zaman bilgisi kullanılamamıştır. MIMIC-IV veya diğer ESK veri setleri içerisinde bu bilginin olması halinde, önerilecek modeller bu ilişkiyi de yakalayabilecek şekilde önerilebilir. Ayrıca zamanla, hastanın durumuna göre kullanılan ilaçların ve klinik notların değişimi analiz edilerek yeni çalışmalar gerçekleştirilebilir.
5. Mortalite ve yoğun bakımda kalma süreleri üzerine eğitilen modellerin tahminleri, hastaları gruplandırarak detaylı bir şekilde incelenebilir. Yaş, cinsiyet, veya geçmiş hastalıklarına göre gruplandırılan hastalar için ayrı ayrı model eğitimi yapılabilir veya önemli özellikler bu gruplara göre çıkartılarak, karsılaşmalar gerçekleştirilebilir. Yapılacak olan bu model eğitimi ve ileri analizlerin, özellikle klinik uzmanlar için önemli ipuçları sağlayacağı düşünülmektedir.
6. Daha gelişmiş derin öğrenme tabanlı mimariler ile (örn: dönüştürücü tabanlı (transformer-based), evrimsel çizge ağları (graph convolutional networks, GNN), dikkat mekanizmalı (attention-based)) deneyler gerçekleştirilerek, deney kümesi genişletilerek yeni çok-kipli model mimarilerinin önerilebileceği düşünülmektedir.

KAYNAKLAR

- [1] Nosowsky,R.andGiordano,T. J.(2006). The Health Insurance Portability and Accountability Act of 1996 (HIPAA) privacy rule: implications for clinical research. In:*Annu. Rev. Med.*57, pp. 575–590.
- [2] Solares,J. R. A.et al. (2020). Deep learning for electronic health records: A comparative review of multiple deep neural architectures. In: *Journal of biomedical informatics*101, p. 103337.
- [3]Shamshirband,S.et al. (2021). A review on deep learning approaches in healthcare systems: Taxonomies, challenges, and open issues. In:*Journal of Biomedical Informatics*113, p. 103627.
- [4]Xiao,C.,Choi,E., andSun,J.(2018). Opportunities and challenges in developing deep learning models using electronic health records data: a systematic review. In:*Journal of the American Medical Informatics Association*25.10, pp. 1419–1428.
- [5]Shickel,B.et al. (2017). Deep EHR: a survey of recent advances in deep learning techniques for electronic health record (EHR) analysis. In:*IEEE journal of biomedical and health informatics*22.5, pp. 1589–1604.
- [6]Shrank,W. H.,Rogstad,T. L., andParekh,N.(2019). Waste in the US health care system: estimated costs and potential for savings. In:*Jama* 322.15, pp. 1501–1509.
- [7]Makary,M. A.andDaniel,M.(2016). Medical error—the third leading cause of death in the US. In:*Bmj*353.
- [8]Davenport,T.andKalakota,R.(2019). The potential for artificial intelligence in healthcare. In:*Future healthcare journal*6.2, p. 94.
- [9] Yu,K. -H.,Beam,A. L., andKohane,I. S.(2018). Artificial intelligence in healthcare. In:*Nature biomedical engineering*2.10, pp. 719–731.
- [10] Choi,E.et al. (2017a). Using recurrent neural network models for early detection of heart failure onset. In:*Journal of the American Medical Informatics Association*24.2, pp. 361–370.
- [11]Awad,A.,Bader–El–Den,M., andMcNicholas,J.(2017). Patient length of stay and mortality prediction: a survey. In:*Health services management research*30.2, pp. 105–120.

- [12] Shang, J. et al. (2019). Gamenet: Graph augmented memory networks for recommending medication combination. In: *proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 33. 01, pp. 1126–1133.
- [13] Dhieb, N. et al. (2020). A secure ai-driven architecture for automated insurance systems: Fraud detection and risk measurement. In: *IEEE Access* 8, pp. 58546–58558.
- [14] Fu, T., Xiao, C., and Sun, J. (2020). Core: Automatic molecule optimization using copy & refine strategy. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 34. 01, pp. 638–645.
- [15] Bardak, B. and Tan, M. (2015). Prediction of influenza outbreaks by integrating Wikipedia article access logs and Googleflu trend data. In: *2015 IEEE 15th International Conference on Bioinformatics and Bioengineering (BIBE)*. IEEE, pp. 1–6.
- [16] Bardak, B. and Tan, M. (2017). Disease outbreak prediction by data integration and multi-task learning. In: *2017 IEEE Conference on Computational Intelligence in Bioinformatics and Computational Biology*. IEEE.
- [17] Shakirov, V., Solovyeva, K., and Dunin-Barkowski, W. (2018). Review of state-of-the-art in deep learning artificial intelligence. In: *Optical memory and neural networks* 27.2, pp. 65–80.
- [18] Choi, E. et al. (2016a). Multi-layer representation learning for medical concepts. In: *proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1495–1504.
- [19] Mikolov, T. et al. (2013). Efficient estimation of word representations in vector space. In: *arXiv preprint arXiv:1301.3781*.
- [20] Pennington, J., Socher, R., and Manning, C. D. (2014). Glove: Global vectors for word representation. In: *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pp. 1532–1543.
- [21] Choi, E. et al. (2018). Mime: Multilevel medical embedding of electronic health records for predictive healthcare. In: *Advances in neural information processing systems* 31.
- [22] Choi, E. et al. (2017b). GRAM: graph-based attention model for healthcare representation learning. In: *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 787–795.
- [23] Miotto, R. et al. (2016). Deep patient: an unsupervised representation to predict the future of patients from the electronic health records. In: *Scientific reports* 6.1, pp. 1–10.

- [24] Medsker,L.andJain,L. C.(1999). Recurrent neural networks: design and applications. CRC press.
- [25] Choi,E.et al. (2016b). Doctor ai: Predicting clinical events via recurrent neural networks. In:*Machine learning for healthcare conference*. PMLR, pp. 301–318.
- [26] Choi,E.et al. (2016c). Retain: An interpretable predictive model for healthcare using reverse time attention mechanism. In:*Advances in neural information processing systems*29.
- [27] Caballero Barajas,K. L.andAkella,R.(2015). Dynamically modeling patient’s health state from electronic medical records: A time series approach. In:*Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 69–78.
- [28] Lipton,Z. C.et al. (2015). Learning to diagnose with LSTM recurrent neural networks. In:*arXiv preprint arXiv:1511.03677*.
- [29] Song,H.et al. (2018). Attend and diagnose: Clinical time series analysis using attention models. In:*Thirty-second AAAI conference on artificial intelligence*.
- [30] Futoma,J.,Morris,J., andLucas,J.(2015). A comparison of models for predicting early hospital readmissions. In:*Journal of biomedical informatics*56, pp. 229–238.
- [31] Moons,E.et al. (2020). A comparison of deep learning methods for ICD coding of clinical records. In:*Applied Sciences*10.15, p. 5262.
- [32] Walonoski,J.et al. (2018). Synthea: An approach, method, and software mechanism for generating synthetic patients and the synthetic electronic health care record. In:*Journal of the American Medical Informatics Association*25.3, pp. 230–238.
- [33] Choi,E.et al. (2017c). Generating multi-label discrete patient records using generative adversarial networks. In:*Machine learning for healthcare conference*. PMLR, pp. 286–305.
- [34] Esteva,A.et al. (2017). Dermatologist-level classification of skin cancer with deep neural networks. In:*nature*542.7639, pp. 115–118.
- [35] Gulshan,V.et al. (2016). Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. In:*Jama*316.22, pp. 2402–2410.
- [36] Joulin,A.et al. (2016). Bag of tricks for efficient text classification. In:*arXiv preprint arXiv:1607.01759*.
- [37] Le,Q.andMikolov,T.(2014). Distributed representations of sentences and documents. In:*International conference on machine learning*. PMLR, pp. 1188–1196.

- [38] Kim, Y. (2014). Convolutional neural networks for sentence classification. In: *arXiv preprint arXiv:1408.5882*.
- [39] Liu, J., Zhang, Z., and Razavian, N. (2018). Deep ehr: Chronic disease prediction using medical notes. In: *Machine Learning for Healthcare Conference*. PMLR, pp. 440–464.
- [40] Si, Y. and Roberts, K. (2019). Deep patient representation of clinical notes via multi-task learning for mortality prediction. In: *AMIA Summits on Translational Science Proceedings 2019*, p. 779.
- [41] Boag, W. et al. (2018). What’s in a note? unpacking predictive value in clinical note representations. In: *AMIA Summits on Translational Science Proceedings 2018*, p. 26.
- [42] Mullenbach, J. et al. (2018). Explainable prediction of medical codes from clinical text. In: *arXiv preprint arXiv:1802.05695*.
- [43] Johnson, A. E. et al. (2016). MIMIC-III, a freely accessible critical care database. In: *Scientific data* 3.1, pp. 1–9.
- [44] Devlin, J. et al. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. In: *arXiv preprint arXiv:1810.04805*.
- [45] Vaswani, A. et al. (2017). Attention is all you need. In: *Advances in neural information processing systems* 30.
- [46] Alsentzer, E. et al. (2019). Publicly available clinical BERT embeddings. In: *arXiv preprint arXiv:1904.03323*.
- [47] Lee, J. et al. (2020). BioBERT: a pre-trained biomedical language representation model for biomedical text mining. In: *Bioinformatics* 36.4, pp. 1234–1240.
- [48] Huang, K., AlTosaar, J., and Ranganath, R. (2019). Clinicalbert: Modeling clinical notes and predicting hospital readmission. In: *arXiv preprint arXiv:1904.05342*.
- [49] Fraser, K. C. et al. (2019). Extracting umls concepts from medical text using general and domain-specific deep learning models. In: *arXiv preprint arXiv:1910.01274*.
- [50] Bhatia, P. et al. (2019). Comprehend medical: a named entity recognition and relationship extraction web service. In: *2019 18th IEEE International Conference On Machine Learning And Applications (ICMLA)*. IEEE, pp. 1844–1851.
- [51] Zhu, H., Paschalidis, I. C., and Tahmasebi, A. (2018). Clinical concept extraction with contextual word embedding. In: *arXiv preprint arXiv:1810.10566*.
- [52] Neumann, M. et al. (2019). ScispaCy: fast and robust models for biomedical natural language processing. In: *arXiv preprint arXiv:1902.07669*.

- [53] Kormilitzin, A. et al. (2021). Med7: a transferable clinical natural language processing model for electronic health records. In: *Artificial Intelligence in Medicine* 118, p. 102086.
- [54] Wu, S. et al. (2020). Deep learning in clinical natural language processing: a methodical review. In: *Journal of the American Medical Informatics Association* 27.3, pp. 457–470.
- [55] Chan, H. S. et al. (2019). Advancing drug discovery via artificial intelligence. In: *Trends in pharmacological sciences* 40.8, pp. 592–604.
- [56] Ryu, J. Y., Kim, H. U., and Lee, S. Y. (2018). Deep learning improves prediction of drug–drug and drug–food interactions. In: *Proceedings of the National Academy of Sciences* 115.18, E4304–E4311.
- [57] Öztürk, H., Özgür, A., and Ozkirimli, E. (2018). DeepDTA: deep drug–target binding affinity prediction. In: *Bioinformatics* 34.17, pp. i821–i829.
- [58] Rogers, D. and Hahn, M. (2010). Extended-connectivity fingerprints. In: *Journal of chemical information and modeling* 50.5, pp. 742–754.
- [59] Durant, J. L. et al. (2002). Reoptimization of MDL keys for use in drug discovery. In: *Journal of chemical information and computer sciences* 42.6, pp. 1273–1280.
- [60] Jaeger, S., Fulle, S., and Turk, S. (2018). Mol2vec: unsupervised machine learning approach with chemical intuition. In: *Journal of chemical information and modeling* 58.1, pp. 27–35.
- [61] Honda, S., Shi, S., and Ueda, H. R. (2019). Smiles transformer: Pre-trained molecular fingerprint for low data drug discovery. In: *arXiv preprint arXiv:1911.04738*.
- [62] Karpathy, A., Joulin, A., and Fei-Fei, L. F. (2014). Deep fragment embeddings for bidirectional image sentence mapping. In: *Advances in neural information processing systems* 27.
- [63] Ilievski, I. and Feng, J. (2017). Multimodal learning and reasoning for visual question answering. In: *Advances in neural information processing systems* 30.
- [64] Mroueh, Y., Marcheret, E., and Goel, V. (2015). Deep multimodal learning for audio-visual speech recognition. In: *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, pp. 2130–2134.
- [65] Khadanga, S. et al. (2019). Using clinical notes with time series data for ICU management. In: *arXiv preprint arXiv:1909.09702*.

- [66] Shukla,S. N.andMarlin,B. M.(2020). Integrating physiological time series and clinical notes with deep learning for improved ICU mortality prediction. In:*arXiv preprint arXiv:2003.11059*.
- [67] Jin,M.et al. (2018). Improving hospital mortality prediction with medical named entities and multimodal learning. In:*arXiv preprint arXiv:1811.12276*.
- [68] Chen,M.(2017). Efficient vector representation for documents through corruption. In:*arXiv preprint arXiv:1707.02377*.
- [69] Hochreiter,S.andSchmidhuber,J.(1997). Long short-term memory. In:*Neural computation*9.8, pp. 1735–1780.
- [70] Bardak,B.et al. (2021). Prediction of mortality and length of stay with deep learning. In:*2021 29th Signal Processing and Communications Applications Conference (SIU)*. IEEE, pp. 1–4.
- [71] Bardak,B.andTan,M.(2021a). Improving clinical outcome predictions using convolution over medical entities with multimodal learning. In:*Artificial Intelligence in Medicine*117, p. 102112.
- [72] Bardak,B.andTan,M.(2021b). Using Clinical Drug Representations for Improving Mortality and Length of Stay Predictions. In:*2021 IEEE Conference on Computational Intelligence in Bioinformatics and Computational Biology (CIBCB)*. IEEE, pp. 1–8.
- [73] Baltrušaitis,T.,Ahuja,C., andMorency,L. -P.(2018). Multimodal machine learning: A survey and taxonomy. In:*IEEE transactions on pattern analysis and machine intelligence*41.2, pp. 423–443.
- [74] Arrieta,A. B.et al. (2020). Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. In:*Information fusion*58, pp. 82–115.
- [75] Lundberg,S. M.andLee,S. -I.(2017). A unified approach to interpreting model predictions. In:*Advances in neural information processing systems* 30.
- [76] Wang,M.et al. (2020a). An explainable machine learning framework for intrusion detection systems. In:*IEEE Access*8, pp. 73127–73141.
- [77] Parsa,A. B.et al. (2020). Toward safer highways, application of XGBoost and SHAP for real-time accident detection and feature analysis. In:*Accident Analysis & Prevention*136, p. 105405.
- [78] Tjoa,E.andGuan,C.(2020). A survey on explainable artificial intelligence (xai): Toward medical xai. In:*IEEE transactions on neural networks and learning systems*32.11, pp. 4793–4813.
- [79] Rodri´guez-Pérez,R.andBajorath,J.(2019). Interpretation of compound activity predictions from complex machine learning models using

local approximations and shapley values. In:*Journal of Medicinal Chemistry*63.16, pp. 8761–8777.

- [80] Rodríguez-Pérez,R.andBajorath,J.(2020). Interpretation of machine learning models using shapley values: application to compound potency and multi-target activity predictions. In:*Journal of computer-aided molecular design*34.10, pp. 1013–1026.
- [81] Rajkomar,A.et al. (2018). Scalable and accurate deep learning with electronic health records. In:*NPJ Digital Medicine*1.1, pp. 1–10.
- [82] Ramsundar,B.et al. (2019). Deep learning for the life sciences: applying deep learning to genomics, microscopy, drug discovery, and more. O'Reilly Media.
- [83] Jiang,F.et al. (2017). Artificial intelligence in healthcare: past, present and future. In:*Stroke and vascular neurology*2.4.
- [84] Johnson,A. E.et al. (2018). The MIMIC Code Repository: enabling reproducibility in critical care research. In:*Journal of the American Medical Informatics Association*25.1, pp. 32–39.
- [85] Purushotham,S.et al. (2018). Benchmarking deep learning models on large healthcare datasets. In:*Journal of biomedical informatics*83, pp. 112–134.
- [86] Harutyunyan,H.et al. (2019). Multitask learning and benchmarking with clinical time series data. In:*Scientific data*6.1, pp. 1–18.
- [87] Wang,S.et al. (2020b). MIMIC-Extract: A data extraction, preprocessing, and representation pipeline for MIMIC-III. In:*Proceedings of the ACM conference on health, inference, and learning*, pp. 222–235.
- [88] Chung,J.et al. (2014). Empirical evaluation of gated recurrent neural networks on sequence modeling. In:*arXiv preprint arXiv:1412.3555*.
- [89] Shewalkar,A.(2019). Performance evaluation of deep neural networks applied to speech recognition: RNN, LSTM and GRU. In:*Journal of Artificial Intelligence and Soft Computing Research*9.4, pp. 235–245.
- [90] Mangal,S.,Modak,R., andJoshi,P.(2019). LSTM based music generation system. In:*arXiv preprint arXiv:1908.01080*.
- [91] Liu,G.andGuo,J.(2019). Bidirectional LSTM with attention mechanism and convolutional layer for text classification. In:*Neurocomputing*337, pp. 325–338.
- [92] Behera,R. K.et al. (2021). Co-LSTM: Convolutional LSTM model for sentiment analysis in social big data. In:*Information Processing & Management*58.1, p. 102435.

- [93] Su,C.et al. (2020). Neural machine translation with Gumbel tree-LSTM based encoder. In:*Journal of Visual Communication and Image Representation*71, p. 102811.
- [94] Zhang,Y.andYang,J.(2018). Chinese NER using lattice LSTM. In:*arXiv preprint arXiv:1805.02023*.
- [95] Reimers,N.andGurevych,I.(2019). Sentence-bert: Sentence embeddings using siamese bert-networks. In:*arXiv preprint arXiv:1908.10084*.
- [96] Rumelhart,D. E.,Hinton,G. E., andWilliams,R. J.(1985).*Learning internal representations by error propagation*. Tech. rep. California Univ San Diego La Jolla Inst for Cognitive Science.
- [97] Hubel,D. H.andWiesel,T. N.(1959). Receptivefields of single neurones in the cat's striate cortex. In:*The Journal of physiology*148.3, p. 574.
- [98] Hubel,D. H.andWiesel,T. N.(1962). Receptivefields, binocular interaction and functional architecture in the cat's visual cortex. In:*The Journal of physiology*160.1, p. 106.
- [99] Hubel,D. H.andWiesel,T. N.(1968). Receptivefields and functional architecture of monkey striate cortex. In:*The Journal of physiology*195.1, pp. 215–243.
- [100] LeCun,Y.et al. (1989). Backpropagation applied to handwritten zip code recognition. In:*Neural computation*1.4, pp. 541–551.
- [101] Krizhevsky,A.,Sutskever,I., andHinton,G. E.(2012). ImageNet Classification with Deep Convolutional Neural Networks. In:*Advances in Neural Information Processing Systems*. Ed. byPereira,F.et al. Vol. 25. Curran Associates, Inc.
- [102] Deng,J.et al. (2009). Imagenet: A large-scale hierarchical image database. In: *2009 IEEE conference on computer vision and pattern recognition*. Ieee, pp. 248–255.
- [103] Simonyan,K.andZisserman,A.(2014). Very deep convolutional networks for large-scale image recognition. In:*arXiv preprint arXiv:1409.1556*.
- [104] Szegedy,C.et al. (2015). Going deeper with convolutions. In:*Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9.
- [105] He,K.et al. (2016). Deep residual learning for image recognition. In:*Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778.
- [106] Zhao,B.et al. (2017). Convolutional neural networks for time series classification. In:*Journal of Systems Engineering and Electronics*28.1, pp. 162–169.

- [107] Kiranyaz, S. et al. (2019). 1-d convolutional neural networks for signal processing applications. In: *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, pp. 8360–8364.
- [108] Uzuner, Ö. et al. (2011). 2010 i2b2/VA challenge on concepts, assertions, and relations in clinical text. In: *Journal of the American Medical Informatics Association* 18.5, pp. 552–556.
- [109] Sun, W., Rumshisky, A., and Uzuner, O. (2013). Annotating temporal information in clinical narratives. In: *Journal of biomedical informatics* 46, S5–S12.
- [110] Romanov, A. and Shivade, C. (2018). Lessons from natural language inference in the clinical domain. In: *arXiv preprint arXiv:1808.06752*.
- [111] Uzuner, Ö., Luo, Y., and Szolovits, P. (2007). Evaluating the state-of-the-art in automatic de-identification. In: *Journal of the American Medical Informatics Association* 14.5, pp. 550–563.
- [112] Stubbs, A., Kotfila, C., and Uzuner, Ö. (2015). Automated systems for the de-identification of longitudinal clinical narratives: Overview of 2014 i2b2/UTHealth shared task Track 1. In: *Journal of biomedical informatics* 58, S11–S19.
- [113] Stubbs, A. and Uzuner, Ö. (2015). Annotating longitudinal clinical narratives for de-identification: The 2014 i2b2/UTHealth corpus. In: *Journal of biomedical informatics* 58, S20–S29.
- [114] Yadav, V. and Bethard, S. (2019). A survey on recent advances in named entity recognition from deep learning models. In: *arXiv preprint arXiv:1910.11470*.
- [115] Morgan, H. L. (1965). The generation of a unique machine description for chemical structures—a technique developed at chemical abstracts service. In: *Journal of Chemical Documentation* 5.2, pp. 107–113.
- [116] Irwin, J. J. et al. (2012). ZINC: a free tool to discover chemistry for biology. In: *Journal of chemical information and modeling* 52.7, pp. 1757–1768.
- [117] Gaulton, A. et al. (2012). ChEMBL: a large-scale bioactivity database for drug discovery. In: *Nucleic acids research* 40.D1, pp. D1100–D1107.
- [118] team, T. pandas development (Feb. 2020). *pandas-dev/pandas: Pandas*. Version latest.
- [119] Harris, C. R. et al. (Sept. 2020). Array programming with NumPy. In: *Nature* 585.7825, pp. 357–362.
- [120] Pedregosa, F. et al. (2011). Scikit-learn: Machine Learning in Python. In: *Journal of Machine Learning Research* 12, pp. 2825–2830.

- [121] Rehurek,R.andSojka,P.(2011). Gensim–python framework for vector space modelling. In:*NLP Centre, Faculty of Informatics, Masaryk University, Brno, Czech Republic*3.2.
- [122] Abadi,M.et al. (2016). Tensorflow: Large-scale machine learning on heterogeneous distributed systems. In:*arXiv preprint arXiv:1603.04467*.
- [123] Chollet,F.et al. (2015). keras.
- [124] Honnibal,M.andMontani,I.(2017). “spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing”. To appear.
- [125] Mulyar,A.et al. (n.d.). TAC SRIE 2018: Extracting Systematic Review Information with MedaCy. In:*Strain*372 (), p. 338.
- [126] Srivastava,N.et al. (2014). Dropout: a simple way to prevent neural networks from overfitting. In:*The journal of machine learning research* 15.1, pp. 1929–1958.
- [127] Nair,V.andHinton,G. E.(2010). Rectified linear units improve restricted boltzmann machines. In:*Icml*.
- [128] Kingma,D. P.andBa,J.(2014). Adam: A method for stochastic optimization. In:*arXiv preprint arXiv:1412.6980*.
- [129] Agarap,A. F.(2018). Deep learning using rectified linear units (relu). In:*arXiv preprint arXiv:1803.08375*.
- [130] Chollet,F.(2015). Keras.<https://github.com/fchollet/keras>.
- [131] Bolton,E. E.et al. (2008). PubChem: integrated platform of small molecules and biological activities. In:*Annual reports in computational chemistry*. Vol. 4. Elsevier, pp. 217–241.
- [132] Gu,J.et al. (2018). Recent advances in convolutional neural networks. In:*Pattern recognition*77, pp. 354–377.
- [133] Ribeiro,M. T.,Singh,S., andGuestrin,C.(2016). "Why Should I Trust You?": Explaining the Predictions of Any Classifier. In:*Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, pp. 1135–1144.
- [134] Liu,E. -q.andZeng,C. -l.(2021). Blood urea nitrogen and in-hospital mortality in critically ill patients with cardiogenic shock: analysis of the mimic-III database. In:*BioMed Research International*2021.
- [135] Park,H. J.et al. (2021). Efficacy of blood urea nitrogen and the neutrophil-to-lymphocyte ratio as predictors of mortality among elderly patients with genitourinary tract infections: a retrospective multicentre study. In:*Journal of Infection and Chemotherapy*27.2, pp. 312–318.

[136] Taylor, B. C., Wilt, T. J., and Welch, H. G. (2011). Impact of diastolic and systolic blood pressure on mortality: implications for the definition of “normal”. In: *Journal of general internal medicine* 26.7, pp. 685–690.



ÖZGEÇMİŞ

Ad-Soyad: Batuhan Bardak

Uyruđu: T.C.

ÖĞRENİM DURUMU:

- Doktora: 2016-2022, TOBB ETÜ, Bilgisayar Mühendisli ği
- Yüksek Lisans: 2015-2018, ODTÜ, Mühendislik Yönetimi
- Yüksek Lisans: 2014-2016, TOBB ETÜ, Bilgisayar Mühendisli ği
- Lisans: 2009-2014, TOBB ETÜ, Bilgisayar Mühendisli ği

MESLEKİ DENEYİM VE ÖDÜLLER:

Yıl	Yer	Görev
2017-Halen	STM	Veri Bilimci
2014-2017	TOBB ETÜ	Tam Burslu Yüksek Lisans Öğrencisi

YABANCI DİL: İngilizce

TEZDEN TÜRETİLEN YAYINLAR, SUNUMLAR VE PATENTLER:

- Bardak, Batuhan, and Mehmet Tan. "Explainable Prediction of Clinical Tasks with a Deep Multimodal Network Using Clinical Drug Representations on MIMIC-III." Under review in IEEE Journal of Biomedical and Health Informatics, 2022
- Bardak, Batuhan, and Mehmet Tan. "Improving clinical outcome predictions using convolution over medical entities with multimodal learning." Artificial Intelligence in Medicine 117 (2021): 102112.
- Bardak, Batuhan. "Prediction of mortality and length of stay with deep learning." In 2021 29th Signal Processing and Communications Applications Conference (SIU), pp. 1-4. IEEE, 2021.
- Bardak, Batuhan, and Mehmet Tan. "Using Clinical Drug Representations for Improving Mortality and Length of Stay Predictions." In 2021 IEEE Conference on Computational Intelligence in Bioinformatics and Computational Biology (CIBCB), pp. 1-8. IEEE, 2021.

DİĞER YAYINLAR, SUNUMLAR VE PATENTLER:

- Yilmaz, Merve Nur, and Batuhan Bardak. "An Explainable Anomaly Detection Benchmark of Gradient Boosting Algorithms for Network Intrusion Detection Systems." 2022 Innovations in Intelligent Systems and Applications Conference (ASYU). IEEE, 2022.
- Ayan, Emre Tolga, Muhammed Said Zengin, Gamze Deniz, HacıAli Duru, and Batuhan Bardak. "Interpretable Cybersecurity Event Detection in Turkish: A Novel Dataset." 2022 Innovations in Intelligent Systems and Applications Conference (ASYU). IEEE, 2022.
- Tanrısever, Ozer, Emre Tolga Ayan, Muhammed Said Zengin, HacıAli Duru, and Batuhan Bardak. "Named Entity Recognition for Defense Industry." In 2022 29th Signal Processing and Communications Applications Conference (SIU), pp. 1-4. IEEE, 2022.
- Ayan, Emre Tolga, Rabia Arslan, Muhammed Said Zengin, HacıAli Duru, Sedat Salman, and Batuhan Bardak. "Turkish Keyphrase Extraction from Web Pages with BERT." In 2021 29th Signal Processing and Communications Applications Conference (SIU), pp. 1-4. IEEE, 2021.
- Özgül, Ozan Fırat*, Batuhan Bardak*, and Mehmet Tan. "A Convolutional Deep Clustering Framework for Gene Expression Time Series." IEEE/ACM Transactions on Computational Biology and Bioinformatics 18, no. 6 (2020): 2198-2207. *co-first authors

- Tan, Mehmet, Ozan Fırat Özgül, Batuhan Bardak, Işık Ekşioğlu, and Suna Sabuncuoğlu. "Drug response prediction by ensemble learning and drug-induced gene expression signatures." *Genomics* 111, no. 5 (2019): 1078-1088.
- Yagcioglu, Semih, Mehmet Saygin Seyfioglu, Begum Citamak, Batuhan Bardak, Seren Guldamlasioglu, Azmi Yuksel, and Emin Islam Tatli. "Detecting cybersecurity events from noisy short text." *NAACL-HLT(1)*.2019.
- Özgül, Ozan Fırat, Batuhan Bardak, and Mehmet Tan. "Predicting drug activity by image encoded gene expression profiles." In *2018 26th Signal Processing and Communications Applications Conference (SIU)*, pp. 1-4. IEEE, 2018.
- Bardak, Batuhan, and M. Fatih Demirci. "Automatic image selection from images with similar contents." In *2017 25th Signal Processing and Communications Applications Conference (SIU)*, pp. 1-4. IEEE, 2017.
- Bardak, Batuhan, and Mehmet Tan. "Disease outbreak prediction by data integration and multi-task learning." In *2017 IEEE Conference on Computational Intelligence in Bioinformatics and Computational Biology (CIBCB)*, pp. 1-7. IEEE, 2017.
- Bardak, Batuhan, and Mehmet Tan. "Prediction of influenza outbreaks by integrating Wikipedia article access logs and Googleflu trend data." In *2015 IEEE 15th International Conference on Bioinformatics and Bioengineering (BIBE)*, pp. 1-6. IEEE, 2015.