

**TOBB EKONOMİ VE TEKNOLOJİ ÜNİVERSİTESİ**  
**FEN BİLİMLERİ ENSTİTÜSÜ**

**TOPLULUK YÖNTEMİ VE İLAÇ İMZALARI KULLANILARAK ANTİ  
KANSER İLAÇLARIN AKTİVİTE TAHMİNİ**

**YÜKSEK LİSANS TEZİ**

**Ertan TOLAN**

**Bilgisayar Mühendisliği Anabilim Dalı**

**Tez Danışmanı: Yrd. Doç. Dr. Mehmet TAN**

**ARALIK 2016**



Fen Bilimleri Enstitüsü Onayı

.....  
**Prof. Dr. Osman EROĞUL**  
Müdür

Bu tezin Yüksek Lisans derecesinin tüm gereksinimlerini sağladığını onaylarım.

.....  
**Doç. Dr. Oğuz ERGİN**  
Anabilimdalı Başkan Vekili

TOBB ETÜ, Fen Bilimleri Enstitüsü'nün 141111036 numaralı Yüksek Lisans öğrencisi **Ertan TOLAN**'ın ilgili yönetmeliklerin belirlediği gerekli tüm şartları yerine getirdikten sonra hazırladığı "**TOPLULUK YÖNTEMİ VE İLAÇ İMZALARI KULLANILARAK ANTI KANSER İLAÇLARIN AKTİVİTE TAHMİNİ**" başlıklı tezi 19.12.2016 tarihinde aşağıda imzaları olan jüri tarafından kabul edilmiştir.

**Tez Danışmanı:** **Yrd. Doç. Dr. Mehmet TAN** .....  
TOBB Ekonomi ve Teknoloji Üniversitesi

**Jüri Üyeleri:** **Doç. Dr. Pınar KARAGÖZ (Başkan)** .....  
Orta Doğu Teknik Üniversitesi

**Doç. Dr. Osman ABUL** .....  
TOBB Ekonomi ve Teknoloji Üniversitesi



## **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.

Ertan TOLAN



## ÖZET

Yüksek Lisans Tezi

### TOPLULUK YÖNTEMİ VE İLAÇ İMZALARI KULLANILARAK ANTİ KANSER İLAÇLARIN AKTİVİTE TAHMİNİ

Ertan TOLAN

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

Tez Danışmanı: Yrd. Doç. Dr. Mehmet TAN

Tarih: ARALIK 2016

Kişiselleştirilmiş kanser tedavisi, kanserin karmaşıklığı da göz önünde bulundurulduğunda, gelişmekte olan bir yaklaşımdır. Kişiselleştirilmiş tedavinin bir parçası olarak, bir ilacın bir hücre hattındaki etkinliği laboratuvar ortamında ölçülür. Ancak, bu deneylerin yapılması çok zordur ve önemli bir maddi kaynak gerektirir. Bu zorlukların üstesinden gelmek için hesaplayıcı yöntemler, sağlanan veri kümeleri ile birlikte bilgisayar ortamında kullanılır.

Bu çalışmada aktivite tahmini problemi bir regresyon problemi olarak ele alınmıştır ve öncelikle her bir ilaç-hücre hattı çiftinin tahmin hatasının azaltılması için üç farklı regresyon modeli birleştirilerek bir topluluk modeli tasarlanmıştır. Temel modeller; gradyan destekli regresyon, çekirdekli bayes çoklu-iş öğrenme ve iz-norm regülarizasyonlu çoklu-iş öğrenme olarak tanımlanmıştır. Oluşturulan modeli değerlendirmek için iki büyük veri kümesi, genomics of drug sensitivity in cancer ve cancer therapeutics response portal, kullanılmıştır. Bu değerlendirmenin sonuçları, topluluk yönteminin tahminlerinin temel modellerin her birinin tek başlarına yaptıkları tahminlerden önemli ölçüde daha iyi olduğunu göstermektedir. Bunun sonucunda orijinal veri kümelerinde görülmeyen ilaç-hücre hattı çiftleri için oluşturulan modelin sitotoksite tahminleri rapor edilmiştir. Araştırmacılar tarafından yapılan canlı içi (in vivo) laboratuvar çalışmaları da, rapor sonuçlarını desteklemektedir.

İlaç aktivitelerini tahmin etmesi için oluşturulan bir diğler model de imza benzerlik tabanlı regülarizasyonlu çoklu-iş öğrenme modelidir. Bu modelin oluşturulması için öncelikle LINCS veri kümesinde bulunan ilaç-hücre hattı deneyleri incelenerek ilaçlar için aktivite imzaları oluşturulmuştur. Oluşturulan bu ilaç imzaları ile ilaçların arasındaki benzerlikler hesaplanarak, ilaçların benzerliklerini gözardı eden modellerden daha güçlü bir tahmin edici model üretilmiştir. Bu model yine genomics of drug sensitivity in cancer ve cancer therapeutics response portal veri kümeleri kullanarak değerlendirilmiştir ve her iki veri kümesinde de karşılaştırılan modellere göre belirgin bir üstünlük sağlandığı gözlenmiştir.

**Anahtar Kelimeler:** Kanser, Aktivite tahmini, Regresyon, Yapay öğrenme, Çoklu-iş öğrenme, Topluluk modeli, İlaç aktivite imzası



## **ABSTRACT**

Master of Science

### **ACTIVITY PREDICTION OF ANTI CANCER DRUGS BY USING ENSEMBLE LEARNING AND DRUGS' SIGNATURES**

Ertan TOLAN

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

Supervisor: Asst. Prof. Mehmet TAN

Date: DECEMBER 2016

Personalized cancer treatment is an ever-evolving approach due to complexity of cancer. As a part of personalized therapy, effectiveness of a drug on a cell line is measured at the laboratory environments. However, these experiments are backbreaking and money consuming. To surmount these difficulties, computational methods are used with the provided data sets.

In the present study, we considered this as a regression problem and firstly designed an ensemble model by combining three different regression models to reduce prediction error for each drug-cell line pair. We defined our base models as gradient boosting regression, kernelized bayesian multi-task learning and trace-norm regularized multi-task learning. Two major data sets, genomics of drug sensitivity in cancer and cancer therapeutics response portal, were used to evaluate our method. Results of this evaluation show that predictions of ensemble method are significantly better than models *per se*.

Furthermore, we report the cytotoxicity predictions of our model for the drug-cell line pairs that do not appear in the original data sets.

The another method to predict anti cancer drug activity is similarity of signature regularized multi-task learning. To constitute this model, firstly, drug signature is generated by examining drug-cell line experiments found in LINCS data set. Then, by using these activity signatures of drugs, similarities between drugs are calculated and a powerful model which overperforms the models which ignore drug similarities is designed.

Also this model is evaluated with genomics drug sensitivity in cancer and cancer therapeutics response portal data sets and results of the both data sets show that constituted model have significantly more predictive power than contrasted model.

**Keywords:** Cancer, Activity prediction, Regression, Machine learning, Multi-task learning, Ensemble learning, Drug activity signature

## TEŐEKKÜR

Yüksek lisans eğitimin ve tez çalışmalarım boyunca beni destekleyen değerli hocam Yrd. Doç. Dr. Mehmet TAN 'a,

Öğrenim hayatım boyunca sağladığı burs imkanı ile ve de sunduğu çalışma ortamıyla beni destekleyen TOBB Ekonomi ve Teknoloji Üniversitesi ailesine,

Tecrübelerinden faydalandığım TOBB Ekonomi ve Teknoloji Üniversitesi bilgisayar mühendisliği bölümünün değerli öğretim üyelerine,

Birlikte çalışmaktan mutluluk duyduğum yüksek lisans arkadaşlarıma sonsuz teşekkürlerimi sunarım.

Bu tez kapsamında yapılan çalışmalar TÜBİTAK (Proje No: 115E274) tarafından desteklenmektedir.



## İÇİNDEKİLER

	<u>Sayfa</u>
<b>ÖZET</b> . . . . .	iv
<b>ABSTRACT</b> . . . . .	vi
<b>TEŞEKKÜR</b> . . . . .	viii
<b>İÇİNDEKİLER</b> . . . . .	ix
<b>ŞEKİL LİSTESİ</b> . . . . .	xi
<b>ÇİZELGE LİSTESİ</b> . . . . .	xii
<b>KISALTMALAR</b> . . . . .	xiii
<b>SEMBOL LİSTESİ</b> . . . . .	xiv
<b>1. GİRİŞ</b> . . . . .	1
<b>2. İLGİLİ ÇALIŞMALAR</b> . . . . .	5
2.1 Kullanılan Veri Türlerine Göre Çalışmalar . . . . .	5
2.2 Modelleme Yöntemlerine Göre Çalışmalar . . . . .	6
<b>3. TOPLULUK YÖNTEMİNİ KULLANARAK AKTİVİTE TAHMİNİ</b> . . . . .	9
3.1 Temel Alınan Modeller . . . . .	9
3.1.1 Gradyan destekli regresyon . . . . .	10
3.1.2 İz-norm regülerizasyonlu çoklu-iş öğrenme . . . . .	11
3.1.3 Çekirdekli bayes çoklu-iş öğrenme . . . . .	12
3.2 Topluluk Modeli . . . . .	15
<b>4. İLAÇ İMZALARINI KULLANARAK AKTİVİTE TAHMİNİ</b> . . . . .	17
4.1 Lasso Çoklu-İş Öğrenme . . . . .	17
4.2 İlaç Aktivite İmzasının Oluşturulması . . . . .	17
4.3 İlaç Etki Benzerliklerinin Hesaplanması . . . . .	19
4.4 İmza Benzerliği Tabanlı Regülerizasyonlu Çoklu-İş Öğrenme . . . . .	20
<b>5. DENEYSEL SONUÇLAR</b> . . . . .	23
5.1 Ayarlar . . . . .	23
5.2 Veri Kümeleri . . . . .	24
5.2.1 Kanserde İlaç Hassasiyet Genomiği . . . . .	25
5.2.2 Kansere Tedavi Tepki Portalı . . . . .	25
5.2.3 Tümüleşik Ağ Tabanlı Hücre İmza Kütüphanesi . . . . .	26
5.3 Veri Ön İşleme . . . . .	26
5.3.1 Öznitelik seçimi . . . . .	26
5.3.2 Standartlaştırma . . . . .	26
5.3.3 Boyutsal küçültme . . . . .	27
5.4 Topluluk Modeli İçin Çapraz Doğrulama . . . . .	27
5.5 Topluluk Modeli İçin Yeni Aktivite Tahminleri . . . . .	33
5.6 İBTRÇÖ Modeli İçin Çapraz Doğrulama . . . . .	34
<b>6. SONUÇ</b> . . . . .	39
<b>KAYNAKLAR</b> . . . . .	41

<b>EKLER</b> . . . . .	45
<b>ÖZGEÇMİŞ</b> . . . . .	51

## ŞEKİL LİSTESİ

	<u>Sayfa</u>
Şekil 3.1: Topluluk öğrenme modelinin gösterimi. Topluluk öğrenme, farklı öğrenme algoritmalarının çeşitli yollarla birleştirilerek daha genel bir modelin oluşturulduğu öğrenme yöntemidir. . . .	9
Şekil 3.2: Karar ağaçlarının birlikte kullanımı. Her bir ağacın tahmini alınarak bu değerlerin ortalaması çıktı olarak verilmiştir. . . . .	10
Şekil 3.3: Tek-iş öğrenme modelinin gösterimi. Tek-iş öğrenmede, her iş bağımsız olarak değerlendirilir ve bağımsız olarak öğrenilir. . . .	11
Şekil 3.4: Çoklu-iş öğrenme modelinin gösterimi. Çoklu-iş öğrenmede, işler arasındaki ilişki değerlendirilerek aynı anda birden fazla iş öğrenilir. . . . .	13
Şekil 3.5: İkili sınıflandırma için çekirdekli bayes çoklu-iş öğrenme (KBMTL) akış şeması.[18] . . . . .	13
Şekil 3.6: Topluluk modelinin oluşturulması aşamasında kullanılan yığıtlı genelleme yönteminin ilk adımının şematik gösterimi. . . . .	14
Şekil 3.7: Yığıtlı Genelleme yönteminin model tahminlerini birleştirme adımının şematik gösterimi. . . . .	16
Şekil 4.1: İlaç aktivite imzalarının kullanılarak benzerlik matrisinin oluşturulması . . . . .	19
Şekil 5.1: Topluluk yöntemi kullanılarak tahmin edilen GDSC ( $IC_{50}$ ) verisi için ilaçların bireysel karşılaştırılması . . . . .	30
Şekil 5.2: Topluluk yöntemi kullanılarak tahmin edilen GDSC ( $AUC$ ) verisi için ilaçların bireysel karşılaştırılması . . . . .	31
Şekil 5.3: Topluluk yöntemi kullanılarak tahmin edilen CTRP ( $AUC$ ) verisi için ilaçların bireysel karşılaştırılması . . . . .	32
Şekil 5.4: GDSC ( $IC_{50}$ ) için ilaçların bireysel karşılaştırılması . . . . .	35
Şekil 5.5: CTRP ( $AUC$ ) için ilaçların bireysel karşılaştırılması . . . . .	36





## ÇİZELGE LİSTESİ

	<u>Sayfa</u>
Çizelge 5.1: Kullanılan veri kümelerinin bazı özellikleri . . . . .	24
Çizelge 5.2: Topluluk Yöntemi için GDSC veri kümesi( $IC_{50}$ ) sonuçları . . .	28
Çizelge 5.3: Topluluk Yöntemi için GDSC veri kümesi( $AUC$ ) sonuçları . . .	28
Çizelge 5.4: Topluluk Yöntemi için CTRP veri kümesi( $AUC$ ) sonuçları . . .	29
Çizelge 5.5: GDSC veri kümesi için eksik değerlerin tahmini . . . . .	33
Çizelge 5.6: CTRP veri kümesi için eksik değerlerin tahmini . . . . .	34
Çizelge 5.7: İmza benzerlik tabanlı regülarizasyonlu çoklu-iş öğrenme için GDSC veri kümesi ( $IC_{50}$ ) sonuçları . . . . .	34
Çizelge 5.8: İmza benzerlik tabanlı regülarizasyonlu çoklu-iş öğrenme için CTRP veri kümesi ( $AUC$ ) sonuçları . . . . .	35
Çizelge 5.9: CTRP veri kümesinde( $AUC$ ) en fazla sayıda hücre hattı ile deneylenen 40 ilaç için sonuçlar . . . . .	35
Çizelge 5.10: CTRP veri kümesinde( $AUC$ ) en fazla sayıda deneye sahip 40 ilaç için sonuçlar . . . . .	36



## KISALTMALAR

<b>CTRP</b>	: Kanser Tedavi Tepki Portalı
<b>GBR</b>	: Gradyan Destekli Regresyon
<b>GDSC</b>	: Kanserde İlaç Hassasiyet Genomiği
<b>İBTRÇÖ</b>	: İmza Benzerliği Tabanlı Regülerizasyonlu Çoklu-İş Öğrenme
<b>WAMSE</b>	: İlaçların Ortalama Karesele Hatalarının Ortalaması
<b>AMSE</b>	: İlaçların Ortalama Karesele Hatalarının Ağırlıklı Ortalaması
<b>LINCS</b>	: Tümlşik Ağ Tabanlı Hücresele İmza Kütüphanesi
<b>KBMTL</b>	: Çekirdekli Bayes Çoklu-İş Öğrenme
<b>RBF</b>	: Radial Basis Function
<b>SRMTL</b>	: Seyrek Yapı Regülerizasyonlu Çoklu-İş Öğrenme
<b>NDPB</b>	: Tahmin Edicinin En İyi Olarak Tahmin Ettiği İlaç Sayısı
<b>TRMTL</b>	: İz-norm Regülerizasyonlu Çoklu-İş Öğrenme



## SEMBOL LİSTESİ

Bu çalışmada kullanılmış olan simgeler açıklamaları ile birlikte aşağıda sunulmuştur.

### **Simgeler Açıklama**

$AUC$	Doz-tepki eğrisi altındaki alan
$CL$	Hücre hattı
$dr$	İlaç
$gen \uparrow$	Yukarı yönlü regüle edilmiş gen
$gen \downarrow$	Aşağı yönlü regüle edilmiş gen
$IC_{50}$	Yarı maksimum durdurucu konsantrasyon değeri
$MSE$	Mean Squared Error
$\rho_1$	Benzerlik tabanlı regülarizasyon parametresi
$\rho_2$	$\ell_1$ norm regülarizasyon parametresi
$\rho_{L2}$	$\ell_2$ norm regülarizasyon parametresi



## 1. GİRİŞ

Kanser, hücrenin genetik yapısındaki bozulmalar sonucu vücudun çeşitli bölgelerindeki kontrol dışı çoğalma ile oluşan hastalıkların genel adıdır. Hastalığa yakalanma ve ölüm oranı bakımından kanser, dünya çapında önde gelen hastalıklardan biri olarak düşünülmektedir. 2016 yılında ABD’de yaklaşık 1.685.210 yeni kanser vakası teşhisi konacağı ve 595.690 kişinin hastalık dolayısıyla öleceği tahmin edilmektedir [34]. Kanser teşhisinde ve tedavisinde; kanserin ilerlemesine ve ilk gözlemlendiği bölgeye göre çok farklı yöntemler kullanılabilir. Bu tedavi ve teşhis yöntemlerinin başlıcaları cerrahi, kemoterapi, ışın tedavisi, immünoterapi olarak sayılabilir. Bu yöntemler tek başlarına kullanılabildiği gibi birlikte de kullanılabilirler. Örneğin; kemoterapi ve ışın tedavisi çoğu zaman birlikte kullanılan yöntemlerdir.

Bilinen yüzden fazla kanser çeşidi vardır ve genellikle kanserin şekillenmeye başladığı yere göre isim alır. Örneğin, deride ve dokuda başlayıp iç organlara yayılan kanser karsinoma olarak adlandırılırken; kemik, kırkırdak, yağ, kas veya kan damarları gibi bağ veya destek dokularında başlayan kanser sarkoma olarak adlandırılır. Kemik iliği gibi kan yapan dokularda başlayıp çok sayıda anormal kan hücresinin üretilmesine neden olan kanser de lösemi adını alır. Bunun yanında vücutta başladığı yere göre de isim alabilmektedir. Örneğin meme kanseri; meme dokusunda, akciğer kanseri; akciğer dokusunda başlayan ve kontrol edilemeyen hücre çoğalmasındır.

Anti kanser ilaçlar kullanarak bu çeşitli kanser hastalıklarını tedavi etmeye çalışan yöntem olarak bilinen kemoterapi, en çok tercih edilen tedavi yöntemidir [23]. Geleneksel kemoterapide, farklı insanlarda oluşan ve birbirine benzediği düşünülen kanser hastalıklarının tedavisinde benzer ilaçlar kullanılmıştır. Ancak daha sonra bu kanserli hücrelerin genetik olarak birbirine benzemediği yapılan araştırmalarda ortaya konulmuştur [2, 10]. Buradan hareketle yapılan çalışmalarda da aynı tip kanser hastalarına aynı tedavilerin uygulanmasının aslında doğru olmadığı gösterilmiştir [10]. Kontrol dışı çoğalan kanserli hücreleri öldürmeyi amaçlayan anti kanser ilaçları, aynı zamanda sağlıklı hücrelere de zarar verebilmektedir. Kemoterapi tedavisinin hasta üzerinde bıraktığı yan etkiler de düşünüldüğünde, tedavi için doğru ilacın bulunması önem arz etmektedir. Dolayısıyla araştırmacılar, tedavide kullanılacak ilacın çeşidine, dozuna ve ne şekilde kullanılacağına karar vermek için laboratuvar çalışmalarına yoğunlaşmışlardır.

Geleneksel tedavi yöntemleri yerine, kişiselleştirilmiş kanser tedavisi ile, kanser teşhisi konulan hastadan alınan tümör hücresinin yapısı laboratuvar ortamında incelenerek hastanın tedavisine yön verilir. İlgili hasta için, doğru zamanda doğru tedavi yöntemini bulmayı amaçlayan kişiselleştirilmiş kanser tedavisinin, geleneksel terapilere göre daha etkili yöntem olduğu artık bilinen bir gerçektir [19].

Kişiselleştirilmiş tedavinin bir parçası olarak, kanser teşhisi konulmuş bir hastadan alınan tümör hücreleri üzerindeki deneyler, bir tümör hücresinin bir kanser ilacı için ne kadar hassas olduğunu gösterir. Belirli bir ilacın belirli bir kanser tedavisinde etkili olup olmayacağına karar verilmesinde, genellikle kanser hücre hattı üzerindeki bu deneysel sonuçlar başlangıç noktası olur. Ancak, potansiyel ilaç adaylarının sayıca fazla olması nedeniyle, geniş kapsamlı kimyasal bileşik-hücre hattı çiftlerinin deneyleri önemli bir maliyet oluşturmaktadır. Bu sorunun üstesinden gelmek için, tümör hücrelerinin uygulanan ilaçlara tepkisini, laboratuvarında deneyler yaparak ölçmek yerine, ilaç tepkilerini tahmin etmek için kullanılan hesaplama modelleri tasarlanmıştır.

Yapay öğrenme modelleri, son zamanlarda büyük ölçekli ilaç tepki veri tabanlarının [1, 28, 30, 33, 39] yayımlanması sayesinde ilaç aktivite tahmininde kullanılır hale gelmiştir. Bu veri tabanları, yüzlerce kanser hücre hattına karşı çok sayıda kimyasal bileşiğin sitotoksitesite deneylerinin sonuçlarını içermektedir. Hücre sitotoksitesite değerleri, ilaçların hücreleri yok etme kabiliyetinin ölçülmesinde kullanılır. Veri tabanlarında bu veriler önışlemeden geçirilerek hesaplama modelleri için anlamlı hale getirilir. Örneğin sitotoksitesite değeri için GDSC (Genomics of Drug Sensitivity in Cancer) veri kümesinde, doz-tepki eğrisinin altında kalan alana ve yarı maksimum durdurucu konsantrasyon değerinin doğal logaritmasına yer verilir.

Anti kanser ilacının tümör hücresi üzerindeki etkinliđi hakkında bilgi veren doz-tepki eğrisi altındaki alan (*AUC*) ve yarı maksimum durdurucu konsantrasyon (*IC*<sub>50</sub>) değerlerinin tahmin edilmesi için oluşturulan modeller açısından bu problem sınıflandırma ya da regresyon problemi olarak değerlendirilebilir. Sınıflandırmada, bir hücrenin ilaca hassas olup olmadığı öngörölmeye çalışılırken, regresyonda, *IC*<sub>50</sub> veya *AUC* cinsinden sitotoksitesitenin tam değerinin tahmin edilmesi amaçlanır. İkili sınıflandırma (hücre ilaca hassastır veya dirençlidir) ile karşılaştırıldığında, regresyon (ilacın sitotoksitesite değerinin tam olarak tahmini) açıkça daha zordur, ancak bununla beraber ilacın tümör hücresini nasıl etkilediđi konusunda çok daha fazla bilgi verir.

Buradaki en önemli hususlardan biri, hücre hatlarının nasıl karakterize edileceğidir. Hücre hatları; gen ifadesi, DNA metilasyonu ve kopya sayısı varyasyonu verileri gibi birkaç farklı veri türü kullanılarak karakterize edilebilir. Bunların arasında; gen ifadesi verileri, en bilgilendirici olarak görülür ve son zamanlarda yapılan 'DREAM challenge' da bunu doğrulamaktadır [8]. Gen ifadesi profili, hücresel fonksiyonların genel bir görüntüsünü oluşturmak için binlerce genin aktivitesinin ölçülmesi ile elde edilir. Bu profiller, aktif olarak bölünen hücreleri ayırt edebilir veya hücrelerin belirli bir tedaviye nasıl tepki gösterdiğini ölçebilir. Bu sebeple çalışmalarımızda hücre hatları, gen ifadesi profilleri kullanılarak gösterimlenmiştir.

Topluluk öğrenme, farklı öğrenme algoritmalarının çeşitli yollarla birleştirilerek daha genel bir modelin oluşturulduđu öğrenme yöntemidir. Topluluk yöntemlerinde, temel öğrenme algoritmalarının her birinden elde edilebilenden daha iyi tahmin edici performans elde etmek için, çoklu öğrenme algoritmaları kullanılır [11].

Bu tez çalışmasının ilk bölümünde ilaç tepkilerinin tam değerini tahmin etmek için üç farklı yöntemi bir araya getiren bir topluluk modeli önerilmektedir. Temel olarak alınan üç model; gradyan destekli regresyon, çekirdekli bayes çoklu-iş öğrenme ve iz-norm regülarizasyonlu çoklu-iş öğrenmedir. Çoklu-iş öğrenmenin amacı, birden



fazla iş için ortaklaşa öğrenerek öğrenme algoritmalarının performansını arttırmaktır. Anti kanser ilaçlar çoklu-iş öğrenme modellerindeki birbirleriyle ilişkili işler olarak ele alınabildiği için [40], ortak öğrenme kanser ilaçlarının aktivite tahminine uygun bir modeldir ve bu alandaki uygulanabilirliği çok yüksektir. Bu sebeple topluluk modeli için iki tanınmış ve açık kaynak olarak paylaşılmış çoklu-iş öğrenme modeli seçilmiştir; KBMTL (Kernelized Bayesian Multi-task Learning) (Çekirdekli Bayes Çoklu-İş Öğrenme) ve TRMTL (Trace-norm Regularized Multi-task Learning) (İz-norm Regülerizasyonlu Çoklu-İş Öğrenme). Ayrıca GBR (Gradient Boosting Regression) (Gradyan Destekli Regresyon), tahmin gücü dikkate alınarak tek-iş öğrenme modelleri arasından seçilmiştir. Bu üç yöntem birleştirilerek her bir ilacın tahmin hata oranının ve tüm ilaçların ağırlıklı hata ortalamasının düşürülmesi amaçlanmıştır. Tahmin modeli oluşturmak için, GDSC (Genomics of Drug Sensitivity in Cancer) [39] ve CTRP (Cancer Therapeutics Response Portal) [28, 30] veri setleri tarafından sağlanan yüzlerce hücre hattı ve ilaç tepki verisi kullanılmaktadır.

Tez çalışmasının ikinci kısmında, SRMTL (seyrek çizge regülerizasyonlu çoklu-iş öğrenme [43, 44] algoritması kullanılarak ilk kısımda belirtilen problemi ilaç benzerlik ilişkilerini ele alarak çözmeye çalışan başka bir model oluşturulmuştur. SRMTL, parametre olarak verilen çoklu işlerin benzerliği çizgesini kullanarak regülerizasyon işlemi yapan bir modeldir. Bu benzerlik ilişkilerinin modele parametre olarak verilmesi, çoklu-iş öğrenme algoritmasının hedeflediği 'işlerin birlikte öğrenilmesi' görevinin daha iyi bir şekilde yapılmasını amaçlamaktadır. Bu modelin ilaç aktivite tahmini için uygun olduğu düşünülerek imza benzerlik tabanlı regülerizasyonlu çoklu-iş öğrenme (İBTRÇÖ) algoritması, ilaçların hücre hatları üzerindeki hassasiyetinin tahmini için geliştirilmiştir. Gerekli olan ilaç benzerliklerinin hesaplanması için öncelikle her bir ilaç için, ilaçların hücre hatları üzerindeki aktiviteleri göz önünde bulundurularak aktivite imzası oluşturulmuştur. Bu işlem için LINC (The Library of Network-Based Cellular Signatures) [13] veri kümesinden yararlanılmıştır. Bu veri kümesi aracılığıyla, ilaçların hücre hatları ile girdikleri tepkime sonrası gen ifadelerini nasıl değiştirdikleri bilgisine ulaşılabilmektedir. Tepkime sonrası genlerin aşağı ya da yukarı yönlü regüle olduğu bilgisi kullanılarak her ilaç için bir aktivite imzası tasarlanmıştır. Dolayısıyla benzerlik ilişkisi, aynı genleri aynı yönlerde regüle eden iki ilacın birbirine benzediği hipotezine dayandırılarak hesaplanmıştır.

İlaçların benzerlik ilişkisini kullanan bu model, GDSC ve CTRP veri kümeleri üzerinde, bu veri kümelerinde bulunan ilaçların aktivite imzaları ayrı ayrı oluşturularak değerlendirilmiştir. Bu değerlendirme sonucunda her iki veri kümesi için de modelin, benzerlik ilişkisinin verilmediği modellerden daha başarılı sonuçlar verdiği gözlenmiştir.

Bu tez çalışması şu şekilde organize edilmiştir. Bölüm 1'de tez çalışmasında ele alınan problem tanımlanmış ve problemin ortaya çıkışına ve çözümüne yönelik genel bilgiler verilmiştir. Bölüm 2'de problem hakkında yapılan benzer çalışmalar gruplandırılarak özetlenmiştir. Problemin çözümüne yönelik, tez çalışmasında kullanılan yöntemler ve oluşturulan topluluk modeli detaylı olarak Bölüm 3'de ele alınmıştır. Bölüm 4'de imza oluşturma ve benzerlik hesaplama adımları detaylandırılarak, imza benzerlik tabanlı regülerizasyonlu çoklu-iş öğrenme modeli anlatılmıştır. Bölüm 5'de tez çalışmasında kullanılan veri kümeleri ayrıntılı olarak

anlatılmıştır. Oluşturulan topluluk modelinin sonuçları temel modeller ile; benzerlik tabanlı regülarizasyonlu modelin sonuçları benzerlik bilgisinin verilmediği model ile karşılaştırmalı olarak verilmiştir. Tez çalışmasının sonuç kısmı ise Bölüm 6'de bildirilmiştir. Bu bölümde ayrıca oluşturulan modellerin deneysel sonuçları incelenerek, potansiyel gelecek çalışmalara yer verilmiştir.

## 2. İLGİLİ ÇALIŞMALAR

Bu bölümde anti kanser ilaçları için aktivite tahmini yapan literatürdeki çalışmalar incelenecektir. Çalışmalar genel olarak kullandıkları veri türleri ve modelleme teknikleri açısından farklılık gösterir.

### 2.1 Kullanılan Veri Türlerine Göre Çalışmalar

Farklı veri türü olarak; tekli nükleotid mutasyonlar, gen kopyalama sayıları ve gen ifadesinden oluşan herkese açık ya da özel veri kümeleri ile birçok model geliştirilmiştir [8, 20]. Tek veya çoklu kaynaklara dayanan farklı model örnekleri, örneğin gen ifadesi verilerinin gen kopya sayısı verileri ile kombinasyonu, yalnızca gen ifadesi verileri ile en yaygın kullanılan kaynak olarak bildirilmiştir. Karşılaştırmalı analizler genellikle gen ifadesi verisinin en güçlü öngörme özelliklerini içerdiğini ve bu entegre modellerin ilaç tepki tahminlerinin doğruluğunu dikkate değer bir şekilde artırabildiğini göstermiştir [8].

Wan [36] çalışmasında Cancer Cell Line Encyclopedia (CCLE) için iki, DREAM-Challenges için beş farklı genomik karakterizasyon veri kümesi kullanarak heterojen bir yapı oluşturmuştur. CCLE veri kümesi için; gen ifadesi verileri ve mutasyon bilgilerini içeren SNP6 verileri kullanılmıştır. Bu iki veri kümesinden, ilaçların insan kinazlarını hedef aldığını göz önünde bulundurarak, yalnızca kinaz üretilen 400 gen öznitelik olarak seçilmiştir. Rastgele orman yöntemiyle CCLE veri kümesi üzerinde yaptığı değerlendirmede, çoklu genomik karakterizasyonun hata oranını düşürdüğünü belirtmiştir. DREAM-Challenges verileri için ise 5 farklı veri kümesinden oluşturulan bütün kombinasyonlardan birer sonuç üreterek, bu sonuçları topluluk yöntemiyle birleştirmiştir. Burada veri kümesi sayısının artırılması, deney maliyetini de arttırmaktadır. Ayrıca ileride verilerin artacağı da düşünüldüğünde tüm veri kümelerinin kombinasyonundan sonuç olarak bunları topluluk modeli olarak kullanmak çalışma zamanı açısından da problem oluşturabilir.

Genelde öznitelik olarak hücre hatlarının genomik bilgileri kullanılırken, hücre hatlarının genomik karakterizasyonunu ilaçların yapısal özellikleri ile birleştirerek oluşturulmuş modeller de vardır[22] . Ek olarak, ilaçların kimyasal bilgilerinin ve hücre hatlarının profil verisinin birleşiminin girdi olarak alınarak ilaç duyarlılığının tahmininin geliştirilebileceği, NCI-60 verileri üzerinde ilaçlar için modeller eğitilerek gösterilmiştir[7] . Başka bir çalışmada da araştırmacılar modelini hücre hatlarının gen ifadesi seviyelerinin çok değişkenli etkileşimiyle ilişkilendirerek ilaçlar için aktivite tahminini güçlendirmiştir[29] .Bu çalışmalarda genel olarak ileride veri sayısı ve çeşidinin artmaya devam edeceği göz önünde bulundurulmuştur. Oluşturulacak olan

modelde tek bir veri yerine birden fazla veri kullanılmasının tahmin gücünü iyileştireceği belirtilmektedir. Bu modellerin ileride, veri miktarı ve çeşidi arttıkça, daha iyi performans gösterecekleri düşünülebilir.

## 2.2 Modelleme Yöntemlerine Göre Çalışmalar

İlaç aktivite tahmini için tasarlanan modeller daha çok gözetici öğrenme tekniklerine dayanmaktadır. Ancak gözetici tekniklerden de, gözetici tahmin modellerinin oluşturulmasında yararlanılabilmektedir [4, 17, 25]. Ayrıca çalışmalar çoklu-iş öğrenme ve tekli-iş öğrenme modelleri olarak da incelenebilir. Çoklu-iş öğrenme modeli, benzer ilaçların öğrenme aşamasında birbirinden yararlanmasını amaçlarken, tekli-iş öğrenmede modeller ayrı olarak öğrenirler.

Gözetici teknikler genel olarak regresyon ve sınıflandırma modelleri olarak iki ana başlık altında incelenebilir. Regresyon modelleri, hesaplama tekniklerini kullanır ve ilaç hassasiyetini  $IC_{50}$  veya  $AUC$  değerleri cinsinden tahmin etmeye çalışır [14, 24]. Sınıflandırma modellerinde ise önceden belirlenen bir değeri göz önünde bulundurarak ilacın hücre hattı üzerindeki aktivitesi hassas ya da dirençli olarak tahmin edilmeye çalışılır [15, 20].

Anti kanser ilaç aktivite tahmini için çoklu-iş ya da tek-iş öğrenme algoritmaları kullanılmıştır. Çoklu-iş öğrenme modelinin ilaç aktivite tahmini için uygun bir model olduğu gösterilmiştir [40]. Bu çalışmada regülarizasyon için kullanılan iz-norm problemi çarpanları dönüşümlü yönlendirme yöntemi (ADMM) ile çözülerek model üretilmiştir. Üretilen bu model farklı sayıdaki ilaçlardan oluşan kümeler ile eğitilmiştir. Alınan sonuçlar üç farklı veri kümesi kullanılarak elastik net ile karşılaştırıldığında daha iyi bir performans gösterdiği gözlenmiştir.

Tan [35] çalışmasında ilaçlar arasındaki ilişkilerin doğrusal olmadığı hipotezini öne sürmüş ve bu ilişkileri ortaya çıkarmak için modelinde doğrusal olmayan çekirdek kullanmıştır. Bu çalışmada oluşturulan çekirdekli model, iz-norm ile kullanılarak benzer yapıları hedefleyen ilaçlar arasındaki ilişkilerden çoklu-iş öğrenme yöntemiyle faydalanmıştır.

Veri tabanlarındaki hücre hatları ve ilaçlar için tüm deneyler yapılamamıştır ve bu sebeple bazı kayıp değerler bulunmaktadır. Bu kayıp değerli örnekler çıkarıldığında ise veri tabanları küçülmektedir. Modeller bu kayıp verileri kullanıp kullanmamasına göre de farklılaşabilir. Gönen [18], meydana gelen bu küçülmenin önüne geçmek için kayıp değerli örneklerin de kullanıldığı çekirdekli bayes çoklu-iş öğrenme modelini tasarlamıştır. Ayrıca kullandığı çekirdek ile deneysel gürültüleri de azaltmıştır. Gönen [18], çalışmaları sonucunda öğrenme modelinin anti kanser ilaçlarının tepkisini tahmin etmek için oldukça kullanılabilir olduğunu göstermiştir ve çalışmalarını açık olarak paylaşmaktadır. Bu modelin fazla sayıdaki parametre listesi uygulanabilirliğini azalttığı söylenebilir ancak yöntem sahibi tarafından sağlanan parametreler ile birlikte bu tez çalışmasında da kullanılmıştır.

Bunların yanı sıra farklı modelleme teknikleri de kullanılmıştır. İlaç aktivite tahmini için, benzer hücre hatlarının ve benzer ilaçların benzer tepkiler vereceği hipotezini temel alarak hücre hattı benzerlik ağı ve ilaç benzerlik ağı oluşturulmuştur [41].

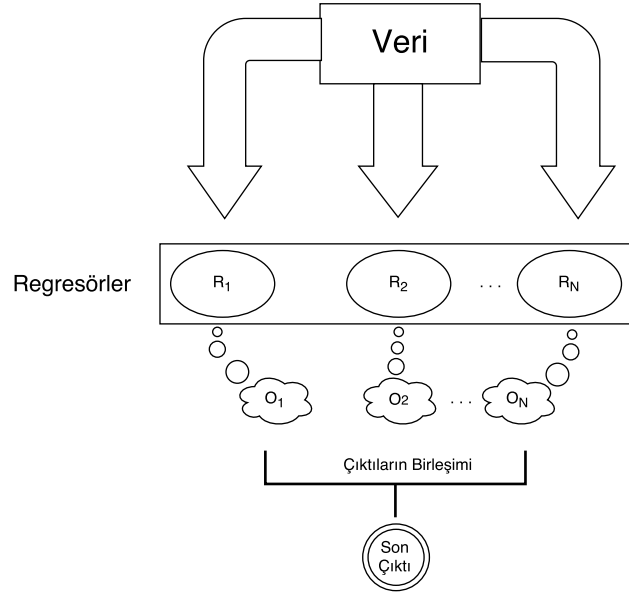
Hücre hattı benzerlik ağını, hücre hatlarının gen ifadesi profillerinin ikili Pearson korelasyonu ile hesaplarken, ilaç benzerlik ağı için ise 1-D ve 2-D ilaç yapılarının ikili Pearson korelasyonu kullanılmıştır. Bu benzerlik ağları doğrusal ağırlıklı bir model ile bütünleştirilerek, tek katmanlı modellerden daha iyi performans gösteren entegre bir ağ önerilmiştir. Farklı ilaçların sinerjik etkilerinden yararlanarak etkili kanser tedavileri geliştirmenin bir başka yöntemi [26] tarafından sunulmuştur. Çalışmalar genel olarak veri tabanları üzerinde ayrı ayrı yapılırken, farklı olarak, Dong [12], kendi modelini oluşturmak ve değerlendirmek için iki farklı veri tabanını, CCLE ve GDSC'yi birleştirmiştir.

Öznitelik seçimi veriyi ilgisiz özniteliklerden arındırarak verinin boyutunu azaltan bir işlemdir. Çoklu örnek öğrenimi (multiple instance learning) ise tek tek örnekler almak yerine, öğrenciye her biri birçok örnek içeren etiketlenmiş bir dizi örnek verir. Zhao, Modelini öznitelik seçimi ve çoklu örnek öğrenimi kullanarak geliştirmiştir[42] .



### 3. TOPLULUK YÖNTEMİNİ KULLANARAK AKTİVİTE TAHMİNİ

Bu bölümde öncelikle kullanılan temel yapay öğrenme yöntemleri ve modellerin oluşturulması aşamasında kullanılan diğer yöntemler ve bunların nasıl seçildikleri, daha sonra ise temel yapay öğrenme algoritmaları birleştirilerek oluşturulan topluluk modeli(Şekil 3.1) açıklanmıştır.



Şekil 3.1: Topluluk öğrenme modelinin gösterimi. Topluluk öğrenme, farklı öğrenme algoritmalarının çeşitli yollarla birleştirilerek daha genel bir modelin oluşturulduğu öğrenme yöntemidir.

#### 3.1 Temel Alınan Modeller

Topluluk modeli oluşturmak için öncelikle üç adet temel model seçilmiştir. Bu modellerin ilki tekli-iş öğrenme (Şekil 3.3) algoritmalarından gradyan destekli regresyondur. Diğer ikisi ise iz-norm regülerizasyonlu ve çekirdekli bayes çoklu-iş öğrenme (Şekil 3.4) algoritmalarıdır.

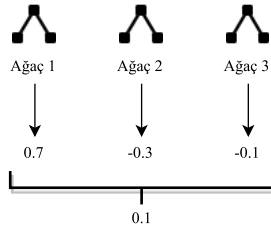
İz-norm regülerizasyonlu çoklu-iş öğrenme yönteminin son zamanlarda kullanıldığı çalışmaların [35, 40] sonuçlarına bakıldığında, bu yöntemde kullanılan regülerizasyonun çoklu-iş öğrenme modelleri için yarar sağladığı görülmüştür. Çekirdekli bayes çoklu-iş öğrenme ise, anti kanser ilaç aktivite tahmini için Gönen [18] tarafından kullanılmıştır ve modelin ortaya koyduğu yenilikler ile iyi bir tahmin edici olduğu gösterilmiştir. Açık kaynak olarak paylaşılan bu iki çoklu-iş öğrenme modeli temel alınan modellerden olmuştur. Gradyan destekli regresyon ise bu

iki çoklu-iş öğrenme modelinin yanında, tekli-iş öğrenme modeli olarak, kendi içinde kullandığı topluluk yönteminden elde ettiği tahmin gücü de göz önünde bulundurulduğunda uygun görülmüştür. Bu yöntemlerin dışında regresyon destek vektör makinesi (SVM regresyon), gauss süreci regresyon (gaussian process regression), lasso regülarizasyonlu çoklu-iş öğrenme gibi yöntemler de denenmiştir. Ancak yapılan çalışmalarda bu üç model uygun görülmüştür.

### 3.1.1 Gradyan destekli regresyon

Regresyon, bir ya da daha fazla bağımsız değişken ile bir bağımlı değişken arasındaki ilişkiyi hesaplamaya yaran analiz işlemidir. Bu hesaplamada doğrusal modeller kullanılabileceği gibi karar ağaçları yardımıyla doğrusal olmayan modeller de üretilebilir. Oluşturulan karar ağaçlarının yapraklarında sürekli değerler yer alır ve iç düğümlerinde ise örnek girdi için verilen öznitelikler ifade edilir. Yapay öğrenmede tek bir karar ağacıyla model oluşturulabilirdiği gibi çok sayıda karar ağacı birlikte kullanılarak topluluk modelleri de oluşturulabilir. Örneğin destek (boosting) algoritmaları zayıf öğrencilerden güçlü bir öğrenci elde etmek için karar ağaçlarını birlikte kullanır (Şekil 3.2).

Gradyan destekli regresyon[16] da çoğunlukla karar ağaçlarıyla temsil edilen ve birçok zayıf öğrencinin bir araya gelmesiyle oluşan bir topluluk modelidir. Zayıf öğrenciler mevcut modele birbiri arkasından eklenerek mevcut öğrencinin eksiklikleri telafi edilir ve zayıf öğrencilerin birleşimiyle güçlü bir öğrenci meydana getirilir.



Şekil 3.2: Karar ağaçlarının birlikte kullanımı. Her bir ağacın tahmini alınarak bu değerlerin ortalaması çıktı olarak verilmiştir.

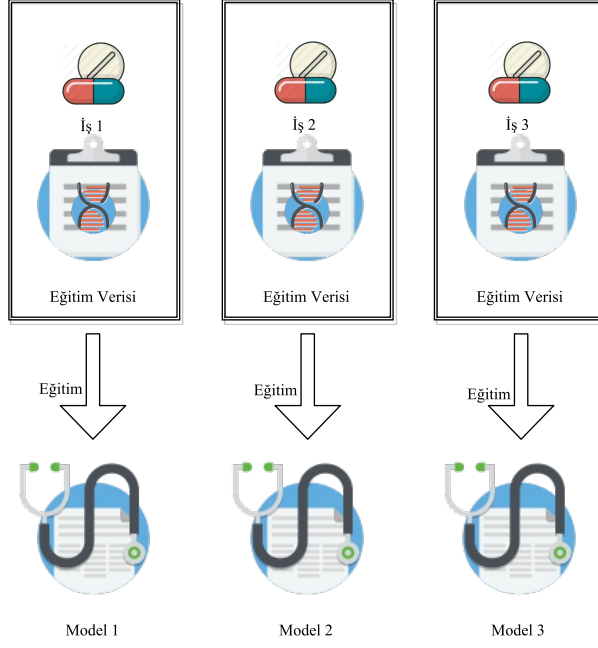
Tez çalışması sırasında oluşturulan topluluk modelinin bir parçası olan gradyan destekli regresyon için MATLAB İstatistik ve Makine Öğrenme Araç Kutusu'nda bulunan LSBoost algoritması kullanılmaktadır. LSBoost algoritmasında *Learners* parametresi ile zayıf öğrencilerin seçimi için esneklik sağlanır. Bu parametre ile zayıf öğrencilerin türü (ağaç, en yakın komşu veya diskriminant) belirlenir. Bu aşamada *Learners* parametresi varsayılan değişkenler ile Ağaç (Tree) olarak ayarlanmıştır.

GBR için diğer iki önemli parametre, sırasıyla modeldeki öğrencilerin sayısı ve büzülme için öğrenme hızını belirten *NLearn* ve *LearnRate*'dir. Gradyan desteği sırasında güncelleştirme için büzülme ile regülarizasyon aşağıdaki gibidir:

$$F_t(x) = F_{t-1}(x) + v \cdot \gamma_t h_t(x), \quad 0 < v \leq 1 \quad (3.1)$$

Burada  $v$  parametresi, öğrenme oranı olarak adlandırılır. Küçük bir öğrenme oranı





Şekil 3.3: Tek-iş öğrenme modelinin gösterimi. Tek-iş öğrenmede, her iş bağımsız olarak değerlendirilir ve bağımsız olarak öğrenilir.

seçmek modelin oluşumu için gelişim sağlar ancak hesaplama zamanı açısından da önemli bir artış olur. *NLearn* ve *LearnRate* arasındaki bu ödünleşimden dolayı bu parametreler birbiriyle ilintili olarak seçilmelidir.

### 3.1.2 İz-norm regülarizasyonlu çoklu-iş öğrenme

Çoklu-iş öğrenme (MTL) (Şekil 3.3), yapay öğrenme modellerinin ayrı ayrı her bir iş için eğitilmesi yerine, birbiriyle ilgili işlerin aynı anda değerlendirilmesidir. Bu şekilde farklı işler için aynı anda eğitmenin daha iyi öğrenmeye yardımcı olduğu [5] tarafından sunulmuştur. Çoklu-iş öğrenme yöntemlerinden faydalanabilmek için işler birbirleriyle ilişki içinde olmalıdır. Antikanser ilaçların modellenmesi de benzer etkiye sahip ilaçların benzer modelleri olabileceği düşüncesiyle çoklu-iş öğrenmeye uygun görülmektedir.

$$\alpha \|w\|_1 = \alpha \sum_i |w_i| \quad (3.2)$$

$$\alpha \|w\|_2^2 = \alpha \sum_i w_i^2 \quad (3.3)$$

Yapay öğrenmede oluşturulan model, eğitim verisine aşırı şekilde uyum sağladığında ve dolayısıyla modelin karmaşıklığının fazla olduğu durumlarda sıkça aşırı öğrenme problemi ile karşılaşılır. Bu problemi aşmak için ise regülarizasyon yönteminin kullanılması uygun bir çözümdür. Regülarizasyon, modelin esnekliğini azaltarak veriye aşırı uyum sağlamayı engelleyen yöntemlerin genel adıdır. Çoklu-iş öğrenme ile kullanılabilen  $\ell_1$ -norm (3.2),  $\ell_2$ -norm (3.3) gibi pek çok regülarizasyon algoritması

vardır. Regülerizasyon kullanan makine öğrenmesi algoritmalarından biri de iz-norm regülerizasyonlu çoklu-iş öğrenmedir. Diğer regülerizasyon fonksiyonları gibi iz-norm regülerizasyon da, öğrenme modellerini, karmaşıklığı cezalandırarak gereğinden fazla uyumu önlemek için ayarlar.

$$\|W\|_p = \left( \sum_{i=1}^{\min\{m,n\}} \sigma_i^p \right)^{1/p} \quad (3.4)$$

İz-norm, aynı zamanda nükleer norm olarak da bilinir ve Schatten p-normunun (3.4)  $p = 1$  olduğu yaygın bir örneğidir. Böylece iz-norm şu şekilde tanımlanabilir:

$$\|W\|_* = \text{iz} \left( \sqrt{W^*W} \right) = \sum_{i=1}^{\min\{m,n\}} \sigma_i \quad (3.5)$$

Burada  $\sigma$ , matrisin köşegeni üzerindeki her bir elaman olarak düşünülebilir. İz-normu temel alan TRMTL [21] ise,  $\lambda$  regülerizasyon parametresi olmak üzere, genel olarak şu eniyileme problemini ele alır:

$$\min_W \sum_{i=1}^n f(W) + \lambda \|W\|_* \quad (3.6)$$

Burada  $f(W)$  yani maliyet fonksiyonunun hesaplanmasında ise en küçük kareler yöntemi kullanılır ve çoklu-iş öğrenme yöntemi için verilen (3.7) problem çözülür.

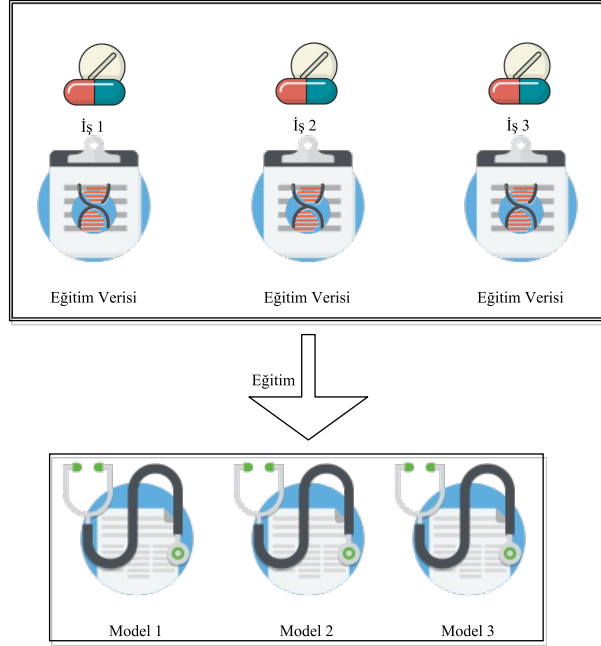
$$\min_W \sum_{i=1}^n \|W_i^T X_i - Y_i\|_F^2 + \lambda \|W\|_* \quad (3.7)$$

MALSAR, Multi-tAsk Learning via StructurAl Regularization, pek çok regülerizasyon algoritmasıyla birlikte çoklu-iş öğrenme yöntemlerinin uygulamasını sağlayan bir araçtır. Tez kapsamında iz-norm regülerizasyonlu çoklu-iş öğrenme metodunu da içeren bu araç kullanılarak modeller oluşturulmuştur.

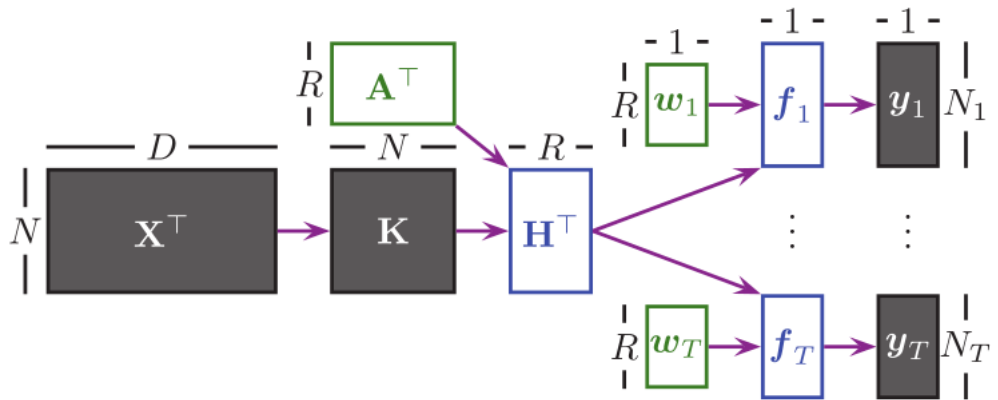
### 3.1.3 Çekirdekli bayes çoklu-iş öğrenme

Çekirdek(kernel) metotları, bir benzerlik fonksiyonu kullanarak veriler arasındaki ilişkileri açığa çıkaran örüntü analizi algoritmasıdır. Yapay öğrenmede çoğu zaman orijinal veri üzerinden sınıflandırma ya da regresyon algoritmalarını kullanmak yerine çekirdek metodu kullanılarak önışlemeden geçirilmiş veri kullanılır. Bu hem verinin boyutunu düşüreceği için daha hızlı çalışmasını sağlar hem de doğrusal olmayan çekirdek kullanıldığında, veriler arasındaki doğrusal olmayan ilişkileri açığa çıkarır.  $\mathbf{X}$  girdi uzayındaki her bir  $\mathbf{x}$  ve  $\mathbf{x}'$  için çekirdek benzerlik fonksiyonu genel olarak şu şekilde verilebilir :

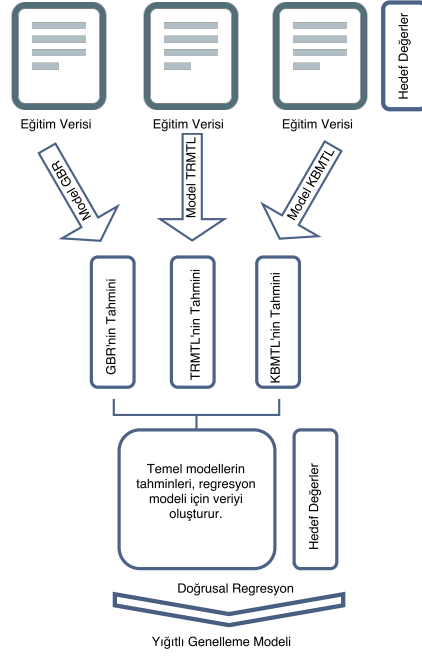
$$K(\mathbf{x}, \mathbf{x}') = \langle \varphi(\mathbf{x}), \varphi(\mathbf{x}') \rangle_{\mathcal{H}} \quad (3.8)$$



Şekil 3.4: Çoklu-iş öğrenme modelinin gösterimi. Çoklu-iş öğrenmede, işler arasındaki ilişki değerlendirilerek aynı anda birden fazla iş öğrenilir.



Şekil 3.5: İkili sınıflandırma için çekirdekli bayes çoklu-iş öğrenme (KBMTL) akış şeması.[18]



Şekil 3.6: Topluluk modelinin oluşturulması aşamasında kullanılan yığıtlı genelleme yönteminin ilk adımının şematik gösterimi.

Burada  $\phi(\mathbf{x})$ , yani  $\mathbf{x}$  in temsilinin nasıl olduğunun açıkça belirtilmesine gerek yoktur ancak,  $\mathbf{x}$  ve  $\mathbf{x}'$  iççarpım için uygun olmalıdır. Ayrıca doğrusal olmayan benzerlik fonksiyonlarıyla birlikte kullanılarak veriler arasındaki doğrusal olmayan ilişkiler de yakalanarak daha verimli algoritmalar tasarlanır. Radyal temelli fonksiyon(3.9), Gönen'in çalışmasında[18] da kullanılan doğrusal olmayan bir benzerlik fonksiyonudur.

$$K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right) \quad (3.9)$$

Radyal temelli fonksiyonun  $\sigma$  parametresinin seçimi önemlidir. Bu parametre veri noktaları arasındaki ikili Öklid uzaklığının ortalaması kullanılarak hesaplanan bir değer ve bu değere komşu diğer dört değer üzerinden yapılan bir iç geçişlemlerle belirlenmiştir.

[18] tarafından verilerde bulunan gürültüyü gidermek ve eksik değerli verileri kullanamama problemini çözmek için her yapay öğrenme işi için ortak bir altuzay kullanımı gibi yenilikleri içeren bir metod önerilmiştir. Çekirdekli bayes çoklu-iş öğrenme (KBMTL) adı verilen bu yöntem hem ikili sınıflandırma (Şekil 3.5) hem de regresyon problemleri için uygundur.

(Şekil 3.5)'de gösterildiği gibi öncelikle  $X$  verisi kullanılarak  $K$  çekirdek matrisi hesaplanır. Daha sonra  $A$  izdüşüm matrisi kullanılarak çekirdek matris bir altuzaya izdüşülür ve  $H$  gizli temsil matrisi bulunur. İkili sınıflandırma kısmında ise gizli temsil üzerinden tahmin yapılır ve bu tahminler sınıf etiketleriyle eşleştirilir.

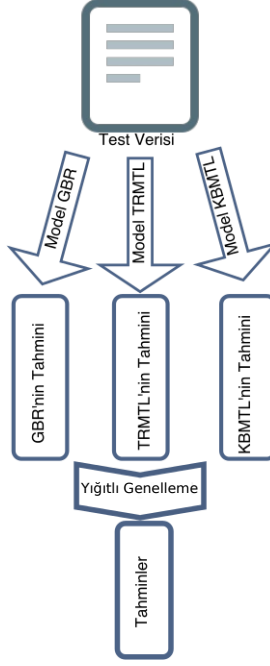
### 3.2 Topluluk Modeli

Toplu öğrenme, daha iyi tahminler elde edebilmek için çeşitli makine öğrenme algoritmalarının birleştirilerek yeni bir model oluşturma yöntemidir. Ortalama alma, oylama, yığıtlama gibi çeşitli topluluk modelleri vardır [31]. Tez çalışmasında yığıtlı genelleme, ortalama alma yöntemi ile karşılaştırılmış ve yığıtlı genellemenin daha iyi sonuçlar verdiği gözlenmiştir. Dolayısıyla bu çalışmada, Bölüm 3.1’de açıklanan temel modeller kullanılarak farklı birleştirme yöntemleri arasından yığıtlı genelleme (stacked generalization) [37] seçilerek topluluk modeli oluşturulmuştur.

Yığıtlı genelleme yönteminde, tahminlerin doğrusal kombinasyonlarının üretilmesi için temel tahmin modellerinin çıktıları ve eğitim verileri için hedef değerler kullanılır. Temel modellerin tahminleri öznitelik vektörü olarak kullanılarak belirlenen katsayıları içeren yığıtlama modeli, daha sonra bu modellerin sonuçlarını birleştirme işleminde kullanılır. Burada katsayısı daha yüksek olan modelin topluluk modelinin genel tahminine etkisinin daha yüksek olması beklenir.

Bu çalışmada yığıtlı genelleme oluşturulurken hem doğrusal regresyon hem de regresyon ağacı denenmiştir ve daha iyi sonuç alınan doğrusal regresyon kullanılmıştır. Ayrıca temel modellerden daha iyi performans almak için eğitim verisi (genelde kullanılan) 2-katlı yerine 5-katlı olarak ele alınmıştır. Burada verinin 5-katlı olarak ele alınması, 5 kere tekrar eden ve her seferinde verinin %80’lik bir kısmı ile modelin eğitilmesini ve kalan %20’lik kısım ile test edilip sonuçların alınmasını ifade eder.

Yığıtlı genelleme, yapay öğrenme modelinde bulunan her iş için ayrı ayrı oluşturulmaktadır. Yani aşağıda verilen adımlar takip edilerek eğitim verisi üzerinden her bir iş için ayrı yığıtlı genelleme modeli oluşturulur. Daha sonra test kısmında her örnek için, temel modellerden alınan tahminler bu yığıtlı genelleme modelinin belirttiği katsayılar kullanılarak hesaplanır.



Şekil 3.7: Yığıtlı Genelleme yönteminin model tahminlerini birleştirme adımının şematik gösterimi.

Temel alınan modelleri uygun şekilde birleştirmek için 5-katlı yığıtlı genelleme, doğrusal regresyon ile kullanıldığında adım adım aşağıdaki gibi ilerler:

- Verinin ilk %80'lik kısmıyla temel modeller oluşturularak %20'lik kısım için hedef değerleri tahmin edilir. Bu işlem 5 kere her bir kat için çalıştırılır.
- Modellerin tahminleri doğrusal regresyon modelinin girdileri, eğitim verisinin hedef değerleri ise bu oluşturulacak regresyon analizi için hedef değerler olarak kullanılır.
- Regresyon analizi sonrası öznitelik olarak ele alınan temel modellerin katsayıları belirlenir. (Şekil 3.6)
- Eğitim verisinin tamamı kullanılarak temel modeller oluşturulur.
- Temel modellerin tahminleri alınır ve belirlenen katsayılar doğrultusunda son tahmin değeri hesaplanır (Şekil 3.7).

İlaç aktivite tahmini için, topluluk yöntemi bu şekilde bir yığıtlı genelleme modeli ile kullanılarak daha güçlü bir tahmin edici oluşturulması amaçlanmıştır. Birbirlerinden farklı üç temel yöntemin yer aldığı bu topluluk modelinde, farklı ilaçlar için yöntemlerin etki katsayıları kendi içlerinde belirlenmiştir.

## 4. İLAÇ İMZALARINI KULLANARAK AKTİVİTE TAHMİNİ

Bu bölümde, öncelikle oluşturulan modelin karşılaştırıldığı Lasso çoklu-iş öğrenme modeline yer verilmiştir. Daha sonra imza benzerlik tabanlı modeli elde etmek için kullanılan ilaç aktivite imzalarının nasıl oluşturulduğu ve bu imzalar kullanılarak hesaplanan benzerlik ilişkisi ele alınmıştır. Son olarak ise benzerlik ilişkisinin seyrek yapı regülarizasyonlu öğrenme modeline eklenerek oluşturulan imza benzerlik tabanlı regülarizasyonlu çoklu-iş öğrenme modeli anlatılmıştır.

### 4.1 Lasso Çoklu-İş Öğrenme

Lasso çoklu-iş öğrenme modeli imza benzerlik tabanlı modelin karşılaştırılması için kullanılan, modelin performans ölçümü için referans alınan yöntemdir. Karşılaştırılan iki model arasındaki tek fark benzerlik ilişkisinin kullanılıp kullanılmaması olmuştur. İmza benzerlik tabanlı regülarizasyonlu çoklu-iş öğrenme modeli, MALSAR tarafından sağlanan Lasso çoklu-iş öğrenme 4.1 modeli ile karşılaştırılarak benzerlik ilişkisinin etkisi ölçülmüştür.

$$\min_W \sum_{i=1}^n \|W_i^T X_i - Y_i\|_F^2 + \rho_2 \|W\|_1 \quad (4.1)$$

Bu modelde adından da anlaşılacağı üzere aşırı öğrenmeyi önleme ve öznelik seçimi için Lasso regülarizasyonu kullanılmıştır.  $\ell_1$  norm olarak da bilinen bu regülarizasyon yöntemi model katsayılarını seyrekleştirmeyi amaçlamaktadır. Bu model ayrıca çoklu-iş öğrenme yöntemi ile Lasso regülarizasyonunu bir arada kullanarak, işler için tüm model katsayılarının eş zamanlı olarak belirlenmesini sağlar.

### 4.2 İlaç Aktivite İmzasının Oluşturulması

LINCS [13] veri tabanında bulunan ilaç-hücre hattı deneylerine, lincscld uygulama programlama arayüzü(API) <sup>1</sup> aracılığıyla erişim sağlanabilmektedir. Bu deneyler, ilaçların hücre hatlarındaki hangi genleri nasıl değiştirdiği bilgisini içerir. Lincscld API' nin, verilen ilaç-hücre hattı çifti için regüle edilmiş probe kümesini sağladığı göz önünde bulundurularak, bir ilacın aktivite imzası oluşturulurken o ilacın veri tabanında bulunan tüm hücre hattı deneyleri sorgulanmıştır.

---

<sup>1</sup><http://api.lincscld.org/>

İlaç etkisiyle değişen probe kümelerine, veri tabanından farklı şekillerde erişilebilir. Tez çalışmasında ilaç aktivite imzası oluşturulurken, bu adımda belirli 1000 probe'tan en çok değişen 50 tanesini içeren küme hesaba katılmıştır. Veri tabanı küratörleri tarafından belirlenen 1000 probe'luk bu liste, gen ifadesindeki bilgiyi genel olarak özetlemektedir<sup>2</sup>. Daha sonra veritabanında probe değişimi cinsinden yer alan bu bilgiler mygene [38] kütüphanesi kullanılarak gen değişimlerine dönüştürülmüştür. Bu noktada *dr* ilaç, *CL* hücre hattı olmak üzere bir (*dr,CL*) çifti için aşağıda verilen şekilde deney imzası (*ExpSig*) oluşturulmuştur:

$$ExpSig(dr, CL) = \langle gen_1 \uparrow, gen_2 \downarrow, gen_3 \downarrow, \dots, gen_K \uparrow \rangle \quad (4.2)$$

Aşağıda örnek olarak doramapimod ilacının MCF7 hücre hattına uygulanması sonucunda; ACAT2 ve ADRB2 genlerinin yukarı yönlü regüle edilmiş olduğu, ADH5 ve ZMIZ1 genlerinin aşağı yönlü regüle edilmiş olduğu gösterilmiştir.

$$ExpSig(doramapimod, MCF7) = \langle ACAT2 \uparrow, ADH5 \downarrow, ADRB2 \uparrow, \dots, ZMIZ1 \downarrow \rangle \quad (4.3)$$

Bir ilaç-hücre hattı çifti için deney parametrelerine bağlı olarak (ilacın ne kadar süre uygulandığı, ilaç dozu vb.) birden fazla deney bulunabilmektedir. Bu farklılaşma ile deney sayısı her çift için farklı sayıda olmaktadır. Örneğin; (*vorinostat, MCF7*) çifti için veri tabanında 188 deney bulunmaktadır. Bu farklılaşmış deneyler tez çalışmasında ayrı deneyler olarak ele alınmıştır ve bu deney imzaları kullanılarak verilen *dr* ilacı için, ilaç aktivite imzası (*ActSig*) oluşturulmuştur.

$$ActSig(dr) = \forall CL \cup_{dr} ExpSig(dr, CL) \quad (4.4)$$

Burada  $\cup_{dr}$ , *dr* ilacı için deney imzalarının birleşimini ifade eder. Bu noktada ilacın gene olan etkisinin keşfedilmesi için kaç kere yukarı, kaç kere aşağı yönlü değişim gösterdiği verisi kullanılmıştır.

$$ActSig(dr) = \langle g_1(\uparrow n_1, \downarrow n_1), \dots, g_K(\uparrow n_K, \downarrow n_K) \rangle \quad (4.5)$$

$g_1(\uparrow n_1)$  ifadesi, *dr* ilacının  $g_1$  genini ( $\uparrow n_1$ ) defa yukarı yönlü, ( $\downarrow n_1$ ) defa aşağı yönlü değiştirdiğini ifade etmektedir. Doramapimod ilacı için örnek bir imza aşağıdaki gibi gösterilebilir:

$$ActSig(doramapimod) = \langle MRPS2(1, 19), MEF2C(12, 4), \dots, SRC(20, 2) \rangle \quad (4.6)$$

Burada, tüm deneylerde zebularine ilacı için, MRPS2 geninin 1 defa aşağı yönlü, 19 defa yukarı yönlü değişime uğradığı ifade edilmektedir.

<sup>2</sup><http://support.lincscloud.org/hc/en-us>



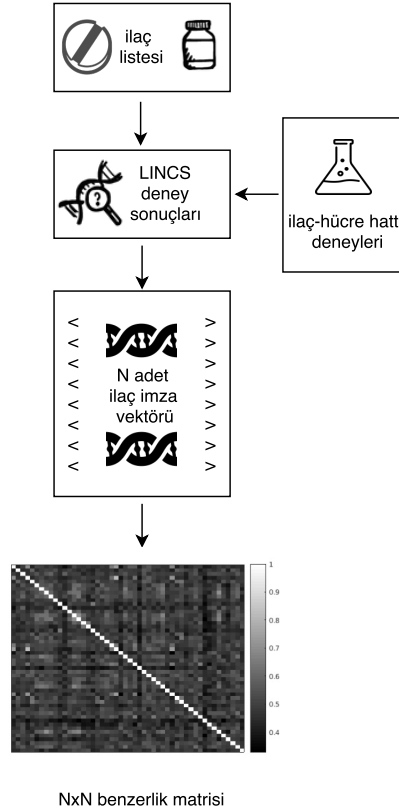
### 4.3 İlaç Etki Benzerliklerinin Hesaplanması

Benzerlik hesaplanması için ilaçların genler üzerindeki etkisinin keşfedildiği aktivite imzaları kullanılırken, benzer ilaçların benzer genleri etkilediği hipotezi esas alınmıştır. Bir ilacın aşağı ya da yukarı yönlü en çok değiştirdiği belirli bir sayıdaki gen listesi ile bir başka ilacın aynı şekilde değiştirdiği gen listesi karşılaştırılarak, bu listelerin kesişimi ne kadar çok ise ilaçlar o kadar benzerdir yorumu yapılmıştır. Bu bilgiler ışığında iki ilacın benzerliği aşağıdaki gibi hesaplanmaktadır:

$$liste_{dr} = \max N(ActSig(dr)) \quad (4.7)$$

$$Benzerlik(dr_A, dr_B) = \frac{|liste_{dr_A} \cap liste_{dr_B}|}{N} \quad (4.8)$$

Burada önemli nokta kesişimleri alınacak gen listeleri için, listelerin uzunluğunun,  $N$ , belirlenmesidir. Bu uzunluk, alınan farklı uzunluklardaki listeler ile yapılan çalışmaların sonuçlarına bakılarak belirlenmiştir. Bu hesaplama tüm ikili ilaç kombinasyonları için yapılarak benzerlik matrisi Şekil 4.1'teki gibi oluşturulmuştur.



Şekil 4.1: İlaç aktivite imzalarının kullanılarak benzerlik matrisinin oluşturulması

Şekil 4.1'de gösterildiği gibi veri kümelerinden sağlanan ilaç listeleri için imza üretmek üzere, ilaçların LINCS üzerinde bulunan deneyleri sorgulanır. Her bir deneyin sonuçlarına bakılarak, o ilacın etkisiyle aşağı ya da yukarı yönde regüle edilmiş genler belirlenir. Belirlenen genler ve regüle edilme sayıları ilacın imzasını

oluşturur. Daha sonra ilaçların benzerlikleri içerdeki ortak gen sayısına bakılarak belirlenir.

#### 4.4 İmza Benzerliği Tabanlı Regülerizasyonlu Çoklu-İş Öğrenme

Çoklu-iş öğrenme metotlarında öğrenilen işlerin birbirleriyle ilişkili olması gerekir. Ancak çoğu durumda verilen işlerin tamamının birbiriyle ilişkili olması beklenemez. Örnek olarak; bazı kanser ilaçları birbirleriyle ilişkili olabilirken, bazı ilaç çiftleri arasında bir ilişki bulunmuyor olabilir. Bu gibi durumlarda da çoklu-iş öğrenme metotlarından iyi bir şekilde yararlanması için, benzerlik ilişkileri girdi olarak alınarak, öğrenilen modellerin katsayılarının regülerizasyonunda bu benzerlikten yararlanır. Bu regülerizasyon yönteminde bağlantılı olduğu düşünülen iş çiftlerinin model katsayıları farkı azaltılmaya çalışılır.

Bu çalışmada kullanılmak üzere MALSAR [44] kütüphanesinde bulunan seyrek çizge regülerizasyonlu çoklu-iş öğrenme modeli seçilmiştir. Bu modelin girdi olarak aldığı çizge yapısı benzerlik ilişkilerini taşıyacak şekilde modele verilmektedir. MALSAR'da bulunan SRMTL metodu; aşağıda verilen çizge yapısı,  $\ell-1$  norm ve  $\ell-2$  norm regülerizasyonlu problemi ele alır:

$$\min_W \sum_{i=1}^n \|W_i^T X_i - Y_i\|_F^2 + \rho_1 \|WR\|_F^2 + \rho_2 \|W\|_1 + \rho_{L2} \|W\|_F^2 \quad (4.9)$$

Burada  $W_i$ ,  $X_i$ ,  $Y_i$  sırasıyla;  $i$  numaralı işin modelini,  $i$  numaralı girdiyi ve  $i$  numaralı işin hedef değerlerini belirtir.  $\rho_1$ ,  $\rho_2$  ve  $\rho_{L2}$  ise model katsayılarının seyrekliğini kontrol eden regülerizasyon parametreleridir.  $\rho_2$  ve  $\rho_{L2}$  parametreleri isteğe bağlı olarak modele verilebilir. Çalışmalarımızda, bu iki parametreden  $\ell-1$  norm için gerekli olan  $\rho_2$  parametresi, çizge seyrekliğini kontrol eden  $\rho_1$  parametresi ile birlikte kullanılarak model oluşturulmuştur.

Problem 4.9'de verilen  $R$  parametresi ise iş benzerlik ilişkilerinin çizge (graph) üzerinde temsil edilmesini sağlar. Bu gösterimde; işlerin her biri birer düğüm (node) olarak düşünülür ve eğer iki iş arasında bir benzerlik varsa, bu iki iş birbirine bir kenar (edge) ile bağlıdır.  $k$  kenar olmak üzere,  $i$  numaralı kenar için,  $A$  ve  $B$  işleri birbirine bağlı ise bu kenar;

$$k_A^{(i)} = \sqrt{\text{Benzerlik}(dr_A, dr_B)} \quad (4.10)$$

ve

$$k_B^{(i)} = (-1) * \sqrt{\text{Benzerlik}(dr_A, dr_B)} \quad (4.11)$$

şeklinde bir vektör ile gösterilmiştir.  $K$  bütün kenar kümesini simgelerse  $R$  çizgesi şu vektörlerden oluşmuştur:

$$R = [k^{(1)}, k^{(2)}, k^{(3)}, \dots, k^{(K)}] \in \mathbb{R}^{t \times \|K\|} \quad (4.12)$$

Bu bilgiler ile birlikte, kenarlardan oluşan  $R$ 'nin kullanıldığı  $\|WR\|_F^2$  ifadesini daha açık bir biçimde yazmak gerekirse :

$$\|WR\|_F^2 = \sum_{i=1}^{\|K\|} \|Wk^{(i)}\|_2^2 = \sum_{i=1}^{\|K\|} \|W_{k_A^{(i)}} - W_{k_B^{(i)}}\|_2^2 \quad (4.13)$$

Dolayısıyla bu regülarizasyon işlemi ile benzer işlerin model kat sayıları arasındaki farkın azaltılması sağlanır. Eğer işler birbirine benzemiyorsa kat sayıları arasındaki fark önemsizdir. Bu durum, R üzerindeki ilgili indeks değerleri sıfır yapılarak sağlanmıştır. İşlerin benzer olduğu durumda ise R üzerindeki ilgili indeks değerleri yüksek değerler olacağı için, algoritma tarafından model katsayıları arasındaki fark küçültülmeye çalışılır.

Regülarizasyon parametresi olan  $\rho_1$ 'in belirlenmesi aşamasında, parametreyi seçmek için grid arama algoritması kullanılmıştır.



## 5. DENEYSEL SONUÇLAR

Öncelikle yapılan deneylerde kullanılan modeller için belirlenen parametreler ve bu parametrelerin nasıl belirlendiği Bölüm 5.1’de ele alınmıştır.

Bölüm 5.2’de anlatılan GDSC ve CTRP veri kümelerinden sağlanan ilaç-hücre hattı aktivite değerleri ve hücre hattı gen ifade profilleri kullanılarak temel alınan modeller ile bu modellerin birleşimiyle oluşturulan topluluk modelinin tahmin güçleri ölçülmüştür. Bu ölçüm öncesi veriyi hazırlamak üzere bazı ön işleme teknikleri kullanılmıştır. Bu teknikler 5.3 bölümünde anlatılmıştır. Daha sonra hazırlanan veri ile oluşturulan modellerin tahmin güçlerinin karşılaştırılması için çapraz doğrulama kullanılmıştır.

Tahmin güçlerine bakıldığında oluşturulan topluluk modelinin ilaç-hücre hattı aktivitesini diğer modellerden daha iyi tahmin ettiği sonucuna varılmıştır. Sonucun ışığında 5.5 bölümde ise veri kümelerinde eksik olarak belirtilen ilaç-hücre hattı çiftleri için tahminler yapılmıştır.

Bölüm 5.6’da ise imza benzerlik tabanlı çoklu-iş öğrenme modeli için çapraz doğrulama sonuçları yer almaktadır. Burada oluşturulan model yine GDSC ve CTRP veri kümeleri kullanılarak test edilmiştir. Oluşturulan modelin performansının karşılaştırılması için referans alınan modelin sonuçları da bu kısımda yer almaktadır.

### 5.1 Ayarlar

Topluluk yönteminin oluşturulması kısmında tamamen MATLAB yazılımı kullanılmıştır. İlaç benzerlik tabanlı yöntemde ise imza oluşturma ve benzerlik hesaplama aşamasında Python programlama dilinden yararlanılmış, modellerin oluşturulması ve değerlendirilmesinde ise yine MATLAB kullanılmıştır.

Topluluk modeli oluştururken kullanılan temel modeller olan, gradyan destekli regresyon, iz-norm regülarizasyonlu çoklu-iş öğrenme ve çekirdekli bayes çoklu-iş öğrenme yöntemleri için model parametreleri ayrı ayrı belirlenmiştir.

Tez kapsamında GBR için, modeldeki öğrenici sayısının belirlendiği *NLearn* ve büzülme için öğrenme hızını belirten *LearnRate* için farklı seçenekler denenmiş ve bu parametreler 100 ve 0.1 olarak ayarlanmıştır . Tanımlanan parametrelerle GBR kullanılarak, iyi bir tahmin modeli elde edilmiştir.

İz-norm için ise regülarizasyon parametresi olan  $\lambda$ 'yı eniyilerken eğitim verisi üzerinde 5 katlı grid arama algoritması kullanılmıştır. [0.1, 1, 10, 100] olarak verilen liste içerisinde parametre seçimi yapılmıştır.

KBMTL ise, çok sayıda parametre isteyen bir yöntem olmasına rağmen, MATLAB uygulamasında verilen, varsayılan parametrelerle kullanılmıştır. Örneğin; alt uzay boyutsallığı 20 olarak ayarlanmıştır.

Veri kümelerinde yer alan ilaçlar için aktivite imzası oluştururken ilaçlara göre deney sayısı farklılık göstermiştir. Sonuç olarak toplam deney sayısının GDSC’de bulunan ilaçlar için 6 - 1057 aralığında, CTRP’de bulunan ilaçlar için 2 - 1000 aralığında değiştiği gözlenmiştir.

Aktivite imzası oluşturma işlemi, GDSC ve CTRP veri kümelerinde bulunan ilaçlar için ayrı ayrı uygulanmıştır. GDSC veri kümesi için 265 ilaçtan 133, CTRP için 481 ilaçtan 286 tanesi için aktivite imzası oluşturulmuştur. Tüm ilaçlar için imza oluşturulamamasının sebebi, LINC5 veritabanında bulunan ilaçlar ile diğer veritabanlarındaki ilaçların eşleştirilememesi ya da LINC5’te ilaçların bulunmamasıdır. Eşleşme problemi, ilaçların farklı veri tabanları için ortak bir kimliğinin bulunmamasından kaynaklanmaktadır. İlaç aktivite imzaları, eşleşen ilaçlar için, benzerlik hesaplanmasında kullanılmak üzere hazır hale getirilmiştir.

Benzerlik hesaplamasında ise önemli nokta, kesişimleri alınacak gen listeleri için listelerin uzunluğunun,  $N$ , belirlenmesidir. Bu uzunluk, alınan farklı uzunluklardaki listeler (10, 20, 40, 80) ile yapılan çalışmaların sonuçlarına bakılmıştır ve 80 olarak belirlenmiştir.

İmza benzerlik tabanlı regülarizasyonlu çoklu-iş öğrenme modelinde kullanılan regülarizasyonların her ikisi ( $\ell_1$  norm ve benzerlik tabanlı) için de eğitim verisi üzerinde yapılan 5 katlı grid arama algoritması kullanılarak parametreler belirlenmiştir. Ayrıca bu yöntemin karşılaştırılmasında kullanılan Lasso çoklu-iş öğrenme modeli için de regülarizasyon parametresi belirlenirken yine grid arama algoritması kullanılmıştır.

## 5.2 Veri Kümeleri

Oluşturulan topluluk modelinin ve imza benzerlik tabanlı modelin tahmin güçlerinin değerlendirilmesi için hücre hattı-ilaç tepkileri ve ilgili hücre hatlarının gen ifadeleri verilerinin sağlandığı iki önemli veri kümesi (GDSC ve CTRP) kullanılmıştır. Kullanılan veri kümelerinin bazı özellikleri karşılaştırmalı olarak Şekil 5.1’de verilmiştir. Son olarak verilen LINC5 veri kümesinden ise ilaç imzalarının oluşturulması aşamasında yararlanılmıştır.

Çizelge 5.1: Kullanılan veri kümelerinin bazı özellikleri

Veri Kümesi	GDSC	CTRP
İlaç Sayısı	265	481
Hücre hattı sayısı	1074	860
Tepki Sayısı	224,510	314,464
$AUC$	✓	✓
$IC_{50}$	✓	✗

### 5.2.1 Kanserde İlaç Hassasiyet Genomiği

Topluluk modelinin ve ilaç benzerlik tabanlı modelin değerlendirilmesi için kullanılan ilk veri kümesi kanserde ilaç hassasiyet genomiği (genomics of drug sensitivity in cancer) (GDSC) [39] kümesidir. GDSC 2012 yılında oluşturulmuştur ve düzenli olarak güncellenmektedir. Bir çok araştırmacı çalışmalarını değerlendirmek için GDSC’de bulunan verileri kullanmaktadır. Tez çalışmasında Temmuz 2016’da yayınlanan güncel sürümü, GDSC v17, kullanılmıştır. Güncellenen sürümünde GDSC, 265 ilaç, 1074 hücre hattı ve 224.510 ilaç tepki değerinden oluşur. İlaçların hücre hatları üzerindeki tepki değerleri hem doz-tepki eğrisi altındaki alan (*AUC*) hem de yarı maksimum durdurucu konsantrasyon değerinin doğal logaritması ( $\log(IC_{50})$ ) cinsinden verildiği için, topluluk modelinin değerlendirme aşamasında her iki hedef değeri de kullanılmıştır. İlaç benzerlik tabanlı model için ise sadece  $\log(IC_{50})$  değeri kullanılmıştır.

Ayrıca, hücre hatlarının gen ifadeleri, RMA normalleştirilmiş bazal ifade profilleri olarak sunulmuştur. Veri kümesinde yer alan hücre hatlarından bazılarının gen ifade verisinde bazı eksikler bulunmaktadır. Gen ifadesi verileri olmayan bu hücre hatları çıkartıldıktan sonra, kalan hücre hattı sayısı 1014’tür. Her ilaç için tüm hücre hatlarının deneyleri yer almadığı için toplam 1014 hücre hattı olmasına rağmen, GDSC’de ilaçlar için deneylenen hücre hattı sayısı ise 363(Rapamycin için) ve 940(Bleomycin-50uM için) aralığında değişkenlik gösterir.

### 5.2.2 Kanser Tedavi Tepki Portalı

Modellerin değerlendirilmesinde kullanılan diğer veri kümesi de kanser tedavi tepki portalıdır. (cancer therapeutics response portal) (CTRP) [28, 30].CTRP’nin ilk sürümü 2012 yılında yayınlanmıştır. Tez çalışmasında kullanılan sürüm ise 2015 yılında yayınlanan ve 2016 yılında bazı önileme işlemleri ile güncellenen sürümüdür. CTRP veri kümesinin güncel sürümünde 481 ilaç ve 860 hücre hattı bulunmaktadır. Bu veri kümesinde ilaç-hücre hattı çiftlerinin duyarlılık skorları sadece doz-tepki eğrisi altındaki alan (*AUC*) cinsinden verilmiştir. Hücre hatları için verilen gen ifadesi değerleri için ise  $\log_2$  dönüşümü yapılmış ortalama değerler kullanılmıştır. Gen ifadesi değerleri, *AUC* ile eşleştirilerek modellerin değerlendirilmesinde kullanılmıştır.

GDSC’dekine benzer olarak CTRP verisinde de bazı hücre hatlarının gen ifadesi verileri eksiktir. Bu hücre hatları veri kümesinden çıkartılarak, 823 hücre hattı üzerinden çalışmalar yapılmıştır. Ayrıca bazı ilaçlar için yetersiz sayıda hücre hattı tepkisi olduğu belirlenmiştir. Bu sebeple belli bir değer altında deneye sahip ilaçların modelde kullanılmaması gerektiği çıkarımıyla yeterli örneğe sahip ilaçlar 439 adet olarak belirlenmiştir. Bu 439 ilaç arasında en az ve en çok deney sonucuna sahip ilaçlar sırasıyla, MG-132(299 deney sonucu) ve Leptomycin(809 deney sonucu)’dir.

### 5.2.3 Tümüleşik Ağ Tabanlı Hücresel İmza Kütüphanesi

Bir hücrenin fenotipinin belirli etmenler tarafından nasıl ve ne zaman değıştirildiğini gözlemek, hastalığa karışan mekanizmalar hakkında ipucu sağlayabilir. Tümüleşik ağ tabanlı hücresel imza kütüphanesi (LINCS) projesi de, bir biyolojik işlemin herhangi birinin bozulmasının, hücrenin moleküler ve hücresel özelliklerinde, davranışında ve işlevinde değışikliklere neden olacağı öncülüne dayanmaktadır. LINCS veri kümesi, kimyasal bileşikler ile tedavi edilen insan hücrelerinin test sonuçlarını içerir. Ayrıca LINCS verileri, bilim insanlarının güncel hastalıklar, ilaçlar ve tedavi yöntemleri ile ilgili problemleri ele almasını kolaylaştırmak için topluluk kaynağı olarak açıkça kullanılabilir hale getirilmiştir. GDSC ve CTRP veri kümelerinde yer alan ilaçlar için imza oluşturulmak amacıyla LINCS veri kümesinde bulunan test sonuçları kullanılmıştır. LINCS veri kümesinde deneylenen çok sayıda bileşik bulunurken toplamda yalnızca 76 hücre hattı ile bu ilaç deneyleri yapılmıştır.

### 5.3 Veri Önışleme

Veri önışleme tekniğı, veri madenciliğı için önemli bir konudur. Normalizasyon, standartlaştırma, gürültü temizleme, öznitelik seçimi gibi teknikler verinin bir sonraki işlem için hazırlanmasını sağlar. Tez çalışmasında aktivite verileri için herhangi bir önışleme yapılmazken, gen ifadeleri için öznitelik seçimi, standartlaştırma ve çekirdek yöntemi kullanılarak boyutsal küçültme yapılmıştır.

#### 5.3.1 Öznitelik seçimi

Gen ifadesi profilleri hücre hatlatı için kullanılan her iki veri kümesi tarafından da ayrı ayrı sağlanmıştır. Bu gen ifadesi verileri çok boyutlu bir veridir, GDSC'den elde edilen veri kümesinde 17.737, CTRP'de ise 18.541 farklı gen için sonuçlar verilmiştir. Bu hem gen sayısı, yani öznitelik sayısı, örnek sayısını aştığı için, hem de bu genlerin tamamının kanser ile ilişkisi olmadığından genler arasında bir seçim yapma ihtiyacı doğmuştur.

Bu ihtiyaç doğrultusunda Malacards [27] veri tabanından yararlanan bir gen seçim prosedürü uygulanmıştır. Bu veritabanını kullanarak, belirli bir hastalıkla ilişkili olduğu bilinen genlerin güncel listesi alınabilir. Buradan yola çıkılarak hücre hatlarının kanser türlerini temel alan bir dizi anahtar kelime üretilmiştir ve kanser hücre hatları ile ilgili 1545 genden oluşan liste indirilmiştir. Tüm modeller için öznitelik olarak, kullanılan veri kümelerinin gen ifadesi verilerindeki gen listesi ile bu listenin kesişiminde yer alan genlerin listesi kullanılmıştır.

#### 5.3.2 Standartlaştırma

Standartlaştırma öznitelikler arasındaki birim farkını ortadan kaldırmak için uygulanan bir tekniktir. Standartlaştırma sonucunda ortalaması sıfır olan ve birim standart sapmaya sahip veri elde edilir. Eğitim verisi için; her bir değerin eğitim



verisinin ortalaması ile farkı alınır ve eğitim verisinin standart sapmasına bölünür. Test verisi için ise; test verisindeki değerler ile eğitim verisinin ortalamasının farkı alınır ve eğitim verisinin standart sapmasına bölünür. Bunun için z skor formülü kullanılmıştır:

$$z = \frac{x - \mu}{\sigma} \quad (5.1)$$

burada  $\mu$  eğitim verisinin ortalaması,  $\sigma$  ise eğitim verisinin standart sapmasıdır.

Bu standartlaştırma işlemi gen ifadesi verileri üzerinde uygulanmıştır. Topluluk modeli, TRMTL ve KBMTL algoritmalarında standartlaştırılmış veriler kullanılmıştır. GBR gibi ağaç tabanlı algoritmalar için ise standartlaştırma işlemine gerek yoktur çünkü bu algoritma sadece verinin daha büyük ya da daha küçük olması ile ilgilenir. İmza benzerlik tabanlı çoklu-iş öğrenme yönteminde ise TRMTL ve KBMTL'de olduğu gibi standartlaştırılmış veriler kullanılmıştır.

### 5.3.3 Boyutsal küçültme

Boyutsal küçültme işlemi verinin öznelik sayısının azaltılmasını sağlayan bir tekniktir. Bu işlem için temel bileşenler analizi yönteminden faydalanılabileceği gibi, çekirdek numarası kullanılarak hem boyutta azaltma sağlanır hem de doğrusal modeller, doğrusal olmayan boyuta taşınır. Bu amaçla radyal temelli fonksiyon (RBF) çekirdeğin kullanılması bilinen bir yöntemdir.

Gönen de çalışmasında [18] yine bu boyutsal küçültme işlemini önermiştir. Buradan yola çıkılarak tez çalışmasının ilk kısmındaki topluluk yöntemi için hem KBMTL hem de TRMTL'de RBF çekirdeği kullanılarak bu boyutsal küçültme işlemi uygulanmıştır. GBR için ise standartlaştırma işleminde olduğu gibi bu teknik de göz ardı edilmiştir. İkinci kısımdaki imza benzerlik tabanlı model için de yine RBF çekirdeği kullanılmıştır.

### 5.4 Topluluk Modeli İçin Çapraz Doğrulama

Topluluk modelinin diğer yöntemlerden daha iyi performans gösterdiğinin doğrulanması için, öğrenme modelleri 10-katlı çapraz doğrulama ile değerlendirilmiştir. Performans karşılaştırması üç farklı ölçev kullanılarak yapılmıştır. Bu ölçevler; ilaçların ortalama karesel hatalarının (5.2) ortalaması (AMSE) (5.3), ilaçların ortalama karesel hatalarının ağırlıklı ortalaması (WAMSE) (5.4) ve her tahmin edicinin en iyi olarak tahmin ettiği ilaç sayısı (NDPB)'dir. AMSE ve WAMSE aşağıdaki gibi hesaplanmıştır.

$$MSE = \frac{1}{n} \sum_i^n (\hat{Y}_i - Y_i)^2 \quad (5.2)$$

$$AMSE = \frac{1}{T} \sum_{t=1}^T \frac{1}{n_t} \sum_{i=1}^n (\hat{Y}_{ti} - Y_{ti})^2 \quad (5.3)$$

Çizelge 5.2: Topluluk Yöntemi için GDSC veri kümesi( $IC_{50}$ ) sonuçları

	AMSE	WAMSE	NDPB
GBR	1.65	1.63	26 / 265
TRMTL	1.87	1.87	4 / 265
KBMTL	1.83	1.81	6 / 265
Topluluk	<b>1.60</b>	<b>1.58</b>	<b>229 / 265</b>

Çizelge 5.3: Topluluk Yöntemi için GDSC veri kümesi( $AUC$ ) sonuçları

	AMSE	WAMSE	NDPB
GBR	$1.51 \times 10^{-2}$	$1.51 \times 10^{-2}$	<b>160</b> / 265
TRMTL	$1.84 \times 10^{-2}$	$1.85 \times 10^{-2}$	18 / 265
KBMTL	$1.72 \times 10^{-2}$	$1.71 \times 10^{-2}$	4 / 265
Topluluk	$1.51 \times 10^{-2}$	$1.5 \times 10^{-2}$	83 / 265

$$WAMSE = \frac{1}{\sum_{t=1}^T n_t} \sum_{t=1}^T \sum_{i=1}^n (\hat{Y}_{ti} - Y_{ti})^2 \quad (5.4)$$

Burada  $T$  ilaç sayısını,  $n$  ise her bir ilacın örnek sayısını belirtir. İlaç hatalarının ortalamasının yanı sıra ağırlıklı ortalamasının verilmesinin sebebi, ilaçların eşit sayıda örnek sayısına sahip olmamasıdır. Yani buradaki ağırlıklar ilaçların deneylendiği hücre hattı sayıdır.

İlaçların aktivite değerlerinin tahmininde hatalar hesaplanırken MSE ölçütü kullanılmıştır. Daha sonra modeller birden fazla ilacın tahmininde kullanıldığı için bu ilaç hatalarının ortalaması alınarak modellerin ilaçların geneli üzerindeki tahmin gücü, modellerin karşılaştırılmasında kullanılmıştır. Burada NDPB ölçütünde ise bir ilacı tahmin eden en iyi modeli bulmak için, o ilaç üzerindeki MSE değeri en küçük olan model seçilmiştir.

GDSC veri kümesinde ilaçların hücre hatları üzerindeki hassasiyeti iki farklı şekilde verilmiştir;  $IC_{50}$  ve  $AUC$ . Dolayısıyla tez çalışmasında da her iki değer üzerinden deneyler yapılmıştır. GDSC verisi için öncelikle hedef olarak belirtilen  $IC_{50}$  değerleri tahmin edilmeye çalışılmıştır. Temel olarak alınan üç modelin ve bunların birleşimiyle oluşturulan topluluk modelinin bu değeri tahmin etmedeki gücü Çizelge 5.2'de gösterilmiştir. Bu çizelgede topluluk modelinin tüm ölçütlerde diğer modellerden daha iyi olduğu görülmektedir. Toplam 265 ilaçtan 229 tanesinde en iyi tahmin performansına sahip model topluluk modeli olmuştur.

Çizelge 5.3'de ise yine GDSC veri kümesi kullanılarak, hedef olarak  $AUC$  değerleri tahmin edilmiştir. Bu tahmin sonuçlarında belirgin bir üstünlükten söz edilemez; topluluk modeli, GBR ile birlikte en iyi performansı göstermiştir.

CTRP veri kümesi ilaçların hücre hatları üzerindeki aktivitesini  $AUC$  değerleri üzerinden paylaşmıştır. Bu sebeple CTRP için deneylerde hedef olarak  $AUC$  değerleri kullanılmıştır. Bu değer tahmini Çizelge 5.4'de modeller arasındaki performansların karşılaştırılması için verilmiştir. Yapılan değerlendirmede topluluk modelinin CTRP

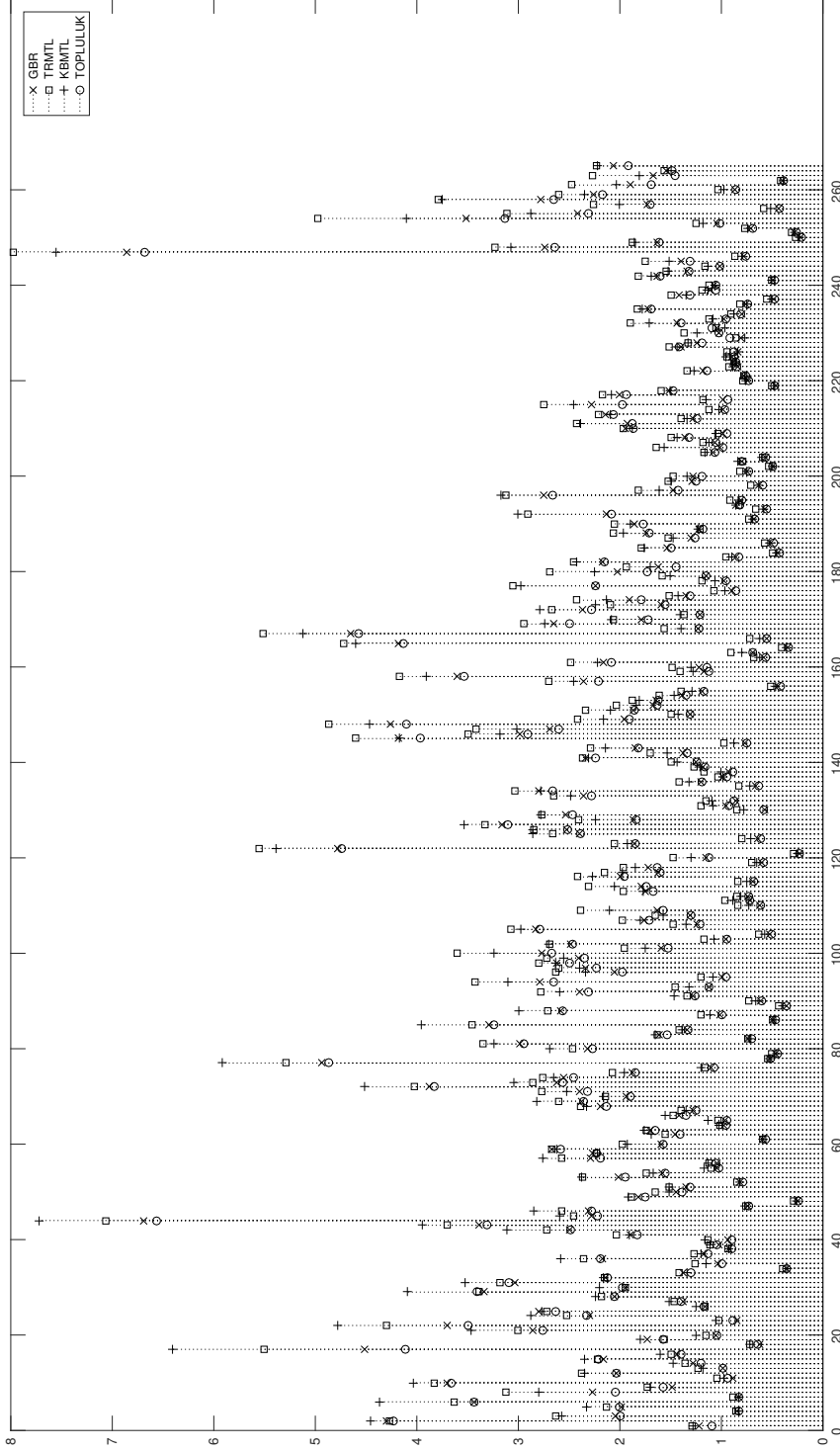
Çizelge 5.4: Topluluk Yöntemi için CTRP veri kümesi(AUC) sonuçları

	<b>AMSE</b>	<b>WAMSE</b>	<b>NDPB</b>
GBR	2.07	2.09	68 / 439
TRMTL	2.38	2.40	63 / 439
KBMTL	2.32	2.33	10 / 439
Topluluk	<b>2.03</b>	<b>2.05</b>	<b>298 / 439</b>

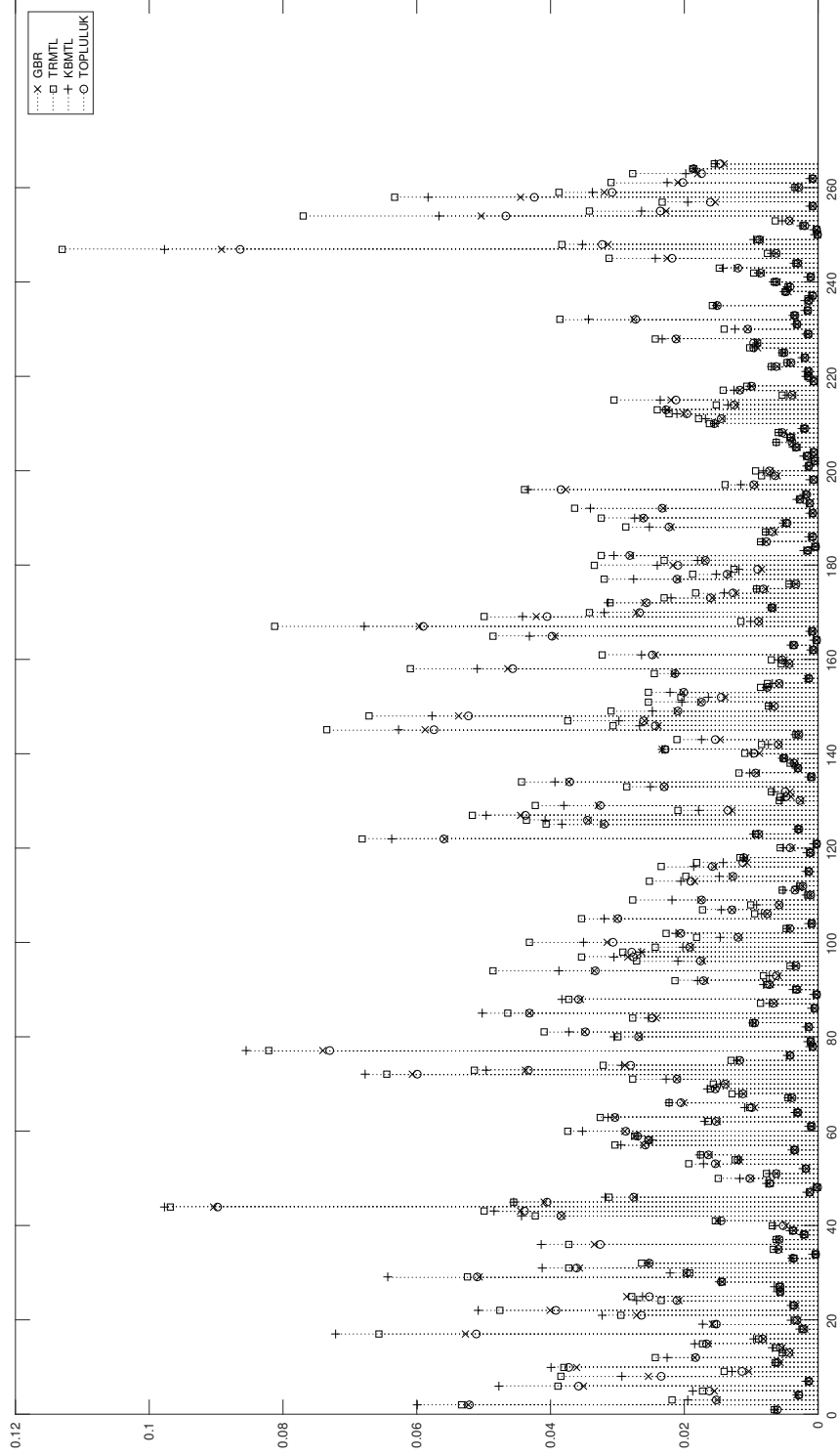
veri kümesi için de olumlu sonuçlar verdiği görülmüştür. Veri kümesinde bulunan 439 ilaçtan 298 tanesi için topluluk modelinin performansının en iyi olduğu gözlenmiştir.

Şekil 5.1 ve 5.2’de GDSC veri kümesinde bulunan 265 ilaç sırasıyla  $IC_{50}$  ve AUC değerlerinin tahmininde hangi modelin hangi ilaç için daha iyi tahminleri ürettiği gösterilmiştir. Buradaki 265 ilacın listesi EK1’de verilmiştir.

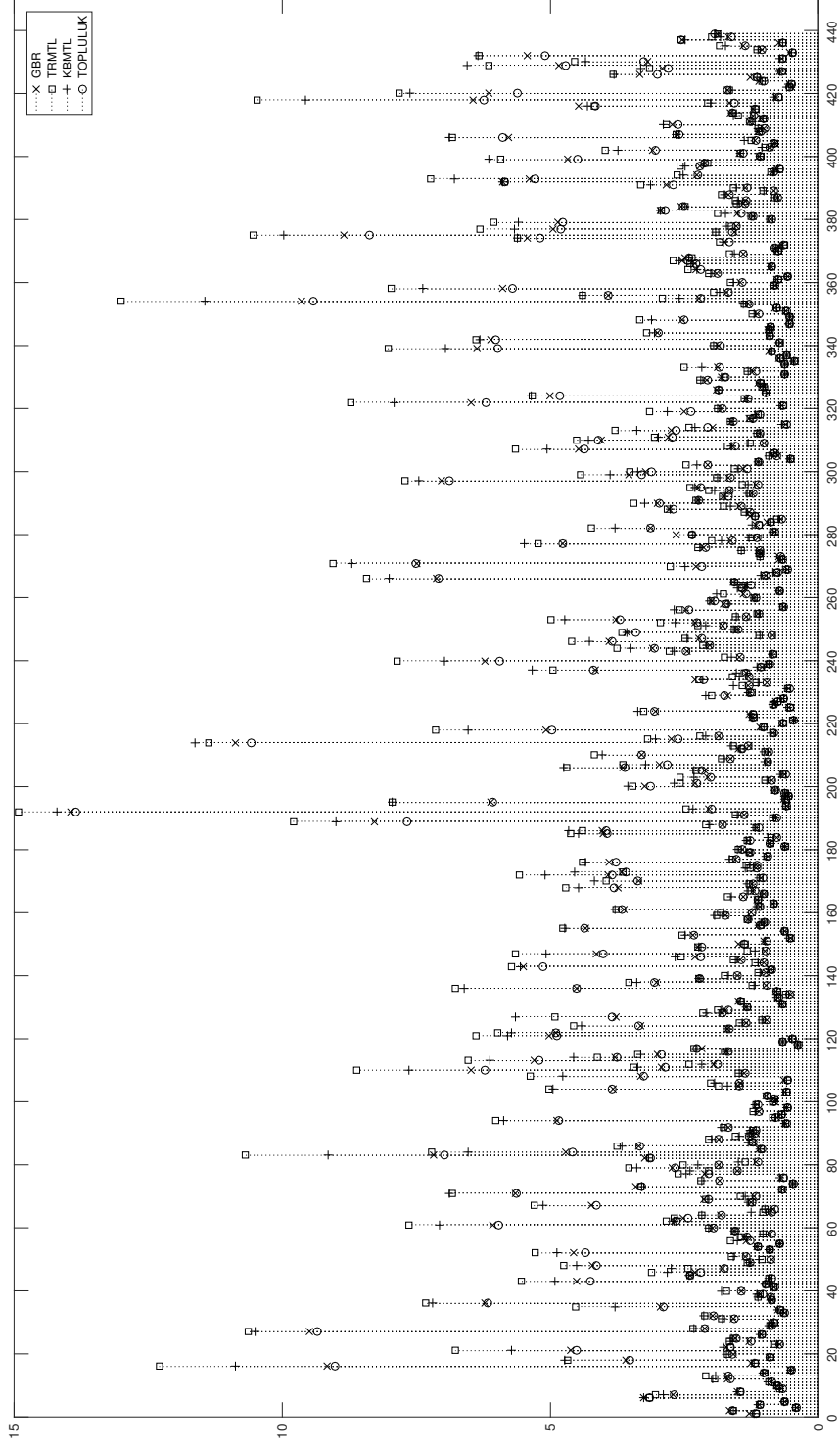
Şekil 5.3 ise CTRP veri kümesinde yer alan 439 ilaç için, modellerin her bir ilacın hücre hatları üzerindeki aktivite tahmini hatasını MSE cinsinden yansıtmaktadır. Buradaki ilaçların listesi de EK2’de verilmiştir.



Şekil 5.1: Topluluk yöntemi kullanılarak tahmin edilen GDSC ( $IC_{50}$ ) verisi için ilaçların bireysel karşılaştırılması



Şekil 5.2: Topluluk yöntemi kullanılarak tahmin edilen GDSC (AUC) verisi için ilaçların bireysel karşılaştırılması



Şekil 5.3: Topluluk yöntemi kullanılarak tahmin edilen CTRP ( $AUC$ ) verisi için ilaçların bireysel karşılaştırılması

Çizelge 5.5: GDSC veri kümesi için eksik değerlerin tahmini

İlaç	Hücre hattı	IC50
[3] Bortezomib	SW756	-7.50
[32]Docetaxel	HSC-3	-6.83
Epothilone B	IOSE-397	-6.07
GSK2126458	RCH-ACV	-5.86
AICAR	A673	-5.68
SN-38	SUP-B15	-5.59
YM155	OCI-LY7	-5.31
Vinorelbine	RCH-ACV	-4.83
Thapsigargin	IOSE-397	-4.79
PD-0325901	G-MEL	-4.62
Ispinesib Mesylate	RCH-ACV	-4.42
Paclitaxel	BICR78	-4.37
Gemcitabine	RCH-ACV	-4.28
Elesclomol	GAMG	-4.23
AP-24534	CML-T1	-4.21
LAQ824	RCH-ACV	-3.84
AUY922	IOSE-397	-3.79
Methotrexate	RS4-11	-3.66
Temsirolimus	SU-DHL-10	-3.65
GW843682X	BICR78	-3.64

### 5.5 Topluluk Modeli İçin Yeni Aktivite Tahminleri

Tüm veri kullanılarak topluluk modeli eğitildikten sonra, ilaçların hücre hatları üzerinde etkili olup olmadığını tanımlanması için, henüz deneylenmemiş ilaç-hücre hattı çiftlerinin değerleri tahmin edilmiştir. Oluşturulan modelin tahmini sonucunda en aktif olacağı düşünülen ilaçlar ve karşılık gelen hücre hatları Çizelge 5.5 ve Çizelge 5.6'de verilmiştir.

Tasarlanan model ile bilgisayar ortamında hangi hücre hatlarının hangi ilaçlara hassas olduğu ortaya çıkartılmıştır. Örneğin SW756 hücre hattı Bortezomib ilacına oldukça hassastır. Ayrıca bir hücre hattı üzerinde hangi ilacın daha aktif olacağı da tasarlanan model ile belirlenebilmektedir. Örnek olarak IOSE-397 hücre hattı üzerinde Epothilone B ilacının Thapsigargin ilacından daha etkili olduğu Çizelge 5.5 'den anlaşılmaktadır.

Uygulanmış canlı içi (in vivo) deneyler, GDSC verisi için oluşturulan modelin Çizelge 5.5'de verilen sonuçlarını desteklemektedir. Örneğin [3] tarafından yapılan bir çalışma, deneylenen tüm ilaçlar arasında Bortezomib ilacının SW756 hücre hattı üzerinde daha etkili olduğunu ortaya çıkarmıştır. Çalışmaya göre; Bortezomib ilacının eeyarestatin ile düşük derişimdeki birleşimi, SW756 kanser hücrelerinin klonal büyümelerini etkili bir biçimde baskı altına almıştır. Aynı zamanda başka bir çalışmada[32] da Docetaxel'in HSC-3 hücrelerindeki iyileştirici etkinliği geliştirmek için ester bağıyla polimer bloğuna konjuge edilmiştir ve Docetaxel ilacının tümörün gelişmesini kontrol ettiği gözlenmiştir.

Çizelge 5.6: CTRP veri kümesi için eksik değerlerin tahmini

İlaç	Hücre hattı	AUC
BRD-K13662825	D283MED	0.05
BRD-K27624156	TE14	1.69
[9] BRD-A05821830	NCIH226	1.78
BRD-A28746609	NCIH226	2.05
BRD-K02130563	AMO1	2.11
BRD-K82109576	NCIH1793	2.50
BRD-K92428232	NCIH1793	3.00
BRD-K23547378	NCIH1793	3.08
BRD-K76674262	NCIH1793	3.16
BRD-K62358710	SKNBE2	3.19
[6] BRD-K15108141	HT29	3.23
BRD-K37764012	D283MED	3.25
BRD-K64890080	NCIH1793	3.32
BRD-K35708212	AMO1	3.746
BRD-K92093830	NCIH1793	3.74
BRD-A35588707	ML1	3.78
BRD-K12343256	L363	3.79
BRD-K34022604	HEC1B	4.17
BRD-K58435339	VMCUB1	4.29
BRD-K29968218	5637	4.40

Çizelge 5.7: İmza benzerlik tabanlı regülarizasyonlu çoklu-iş öğrenme için GDSC veri kümesi ( $IC_{50}$ ) sonuçları

	AMSE	WAMSE	NDPB
Lasso	1.97	1.94	30 / 133
İBTRÇÖ	<b>1.89</b>	<b>1.87</b>	<b>103 / 133</b>

CTRP verisi kullanarak oluşturulan tahminler de benzer şekilde literatürdeki çalışmalar tarafından desteklenmektedir. Örneğin [9], LOR-253 ilacının BRD-A05821830 (Docetaxel) veya BRD-A28746609 (Paclitaxel) ile birlikte NCIH-256 üzerindeki sıralı ve eşzamanlı deneylerinde sinerji gözlenmiştir. Ayrıca [6] çalışmasında da BRD-K15108141 (gemcitabine) ilacının HT29 hücre hattı üzerinde etkili olduğu onaylanmaktadır.

## 5.6 İBTRÇÖ Modeli İçin Çapraz Doğrulama

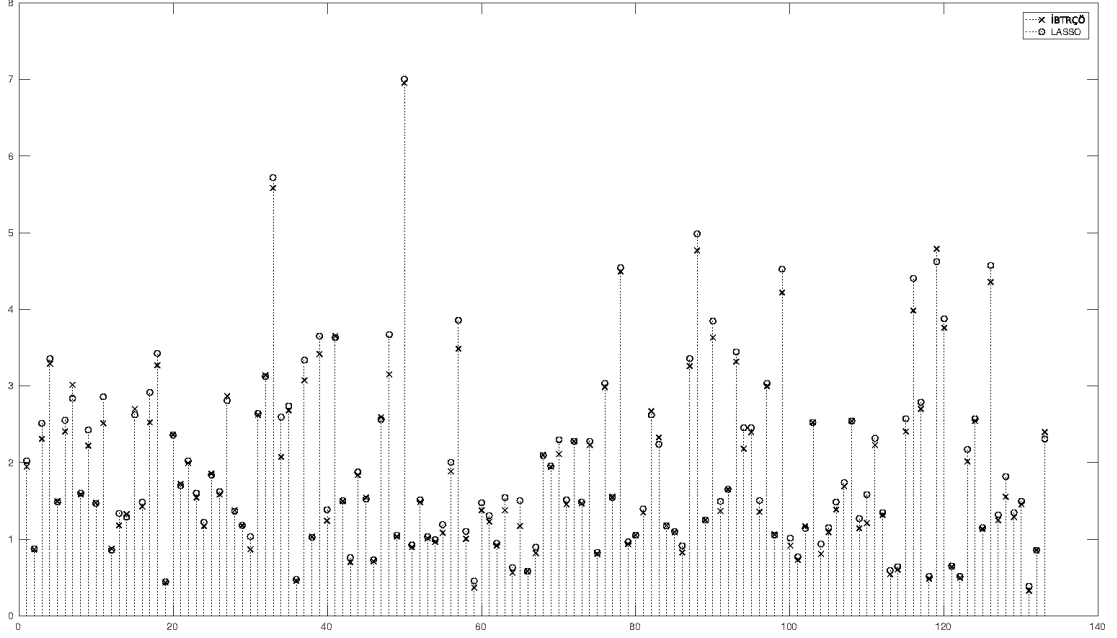
Oluşturulan modelin performansının ölçülmesi için, İBTRÇÖ modeli imzalar arasındaki benzerliğin göz ardı edildiği başka bir çoklu-iş öğrenme modeli (Lasso) ile karşılaştırılmıştır. Topluluk modelinde olduğu gibi üç farklı ölçek kullanılarak yapılan bu karşılaştırma sonucunda imza benzerlik tabanlı regülarizasyonlu çoklu-iş öğrenme modelinin kıyasla daha iyi performans gösterdiği her iki veri kümesi üzerinde de gösterilmiştir (Çizelge 5.7 ve Çizelge 5.8).

İlaçların ortalama karesel hatalarının (5.2) ortalaması (AMSE) (5.3), ilaçların ortalama



Çizelge 5.8: İmza benzerlik tabanlı regülarizasyonlu çoklu-iş öğrenme için CTRP veri kümesi (*AUC*) sonuçları

	AMSE	WAMSE	NDPB
Lasso	2.49	2.54	135 / 286
İBTRÇÖ	<b>2.47</b>	<b>2.53</b>	<b>151 / 286</b>



Şekil 5.4: GDSC ( $IC_{50}$ ) için ilaçların bireysel karşılaştırılması

karesel hatalarının ağırlıklı ortalaması (WAMSE)(5.4) ve her tahmin edicinin en iyi olarak tahmin ettiği ilaç sayısı (NDPB) kullanarak yapılan karşılaştırmada GDSC verisi için daha açık bir üstünlükten söz edilebilir. CTRP için ise hesaplanan sonuçlarda GDSC verisindeki kadar belirgin bir üstünlük görülme de, bu veri kümesi için de sonuçlar olumludur.

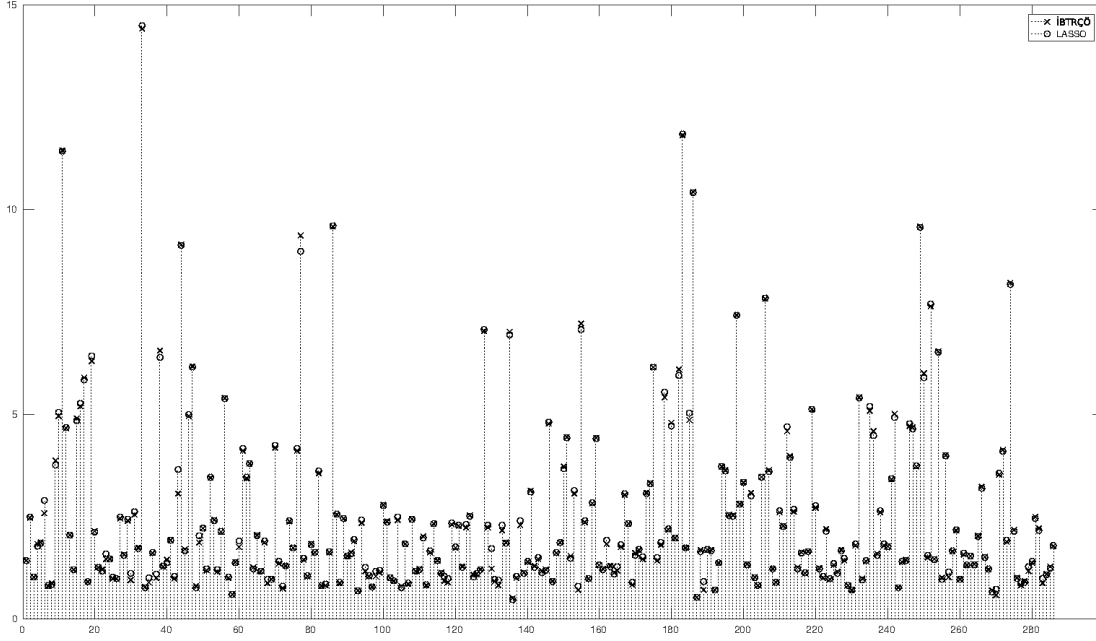
GDSC veri kümesindeki 133 ilaç için ayrı ayrı (Şekil 5.4), CTRP veri kümesindeki 286 ilaç için ayrı ayrı (Şekil 5.5) ilaçların MSE değerlerinin karşılaştırıldığı tablolar ile de bireysel olarak ilaçların tahmin hataları gösterilmiştir.

Ayrıca CTRP veri kümesinde daha az sayıdaki ilaçlar için de bu deneyler tekrarlanmıştır. 286 ilaç içerisinde; en fazla sayıda hücre hattı ile deneylenen 40 ilaç (Çizelge 5.9), en fazla sayıda deneye sahip 40 ilaç (Çizelge 5.10) ile ayrı ayrı yapılan deneylerdeki sonuçlar ile yine oluşturulan modelin katkısı gösterilmiştir.

CTRP verisi için listede bulunan 481 ilaçtan 286'sı için imza oluşturulmuştur. LINC

Çizelge 5.9: CTRP veri kümesinde(*AUC*) en fazla sayıda hücre hattı ile deneylenen 40 ilaç için sonuçlar

	AMSE	WAMSE	NDPB
Lasso	2.62	2.61	14 / 40
İBTRÇÖ	<b>2.59</b>	<b>2.59</b>	<b>26 / 40</b>



Şekil 5.5: CTRP (AUC) için ilaçların bireysel karşılaştırılması

Çizelge 5.10: CTRP veri kümesinde(AUC) en fazla sayıda deneye sahip 40 ilaç için sonuçlar

	AMSE	WAMSE	NDPB
Lasso	3.13	3.18	20 / 40
İBTRÇÖ	<b>3.11</b>	3.18	20 / 40

veri kümesine bakıldığında bu ilaçlardan bazılarının çok az sayıda hücre hattı ile deneylendiği gözlenmiştir. Bunun için, imzaların daha fazla deneylenmiş olmasının daha sağlam imza oluşturulması için gerekli bilgiyi vereceği düşüncesiyle, deney sayısı olarak ve deneylenen hücre sayısı olarak belli bir değerin üstündeki ilaçlar çalışmada kullanılmıştır. Bu çalışmanın sonucunda özellikle deneylenen hücre sayısına göre belirlenen ilaçlarda sonuçların olumlu olduğu gözlenmiştir. Bu da ileride LINCS veri kümesindeki deney sayısı arttığında tüm ilaçlar için daha sağlam modeller oluşturulabileceğinin göstergesidir.



## 6. SONUÇ

Kanser tedavisi için paylaşılan veri kümelerindeki gelişmenin etkisiyle, ilaç hassasiyet analizinde bilgisayar üzerinde yapılan (*in silico*) modellerin kullanımı önem kazanmıştır. Bu durum hem analizlerin daha hızlı bir şekilde yapılmasını sağlamaktadır hem de araştırmacıları ve hastaları oldukça yüksek bir maliyetten kurtarmaktadır.

Bu çalışmanın ilk kısmında ilaç tepkilerini tahmin etmek için üç farklı öğrenme modelini; gradyan destekli regresyon, iz-norm regülarizasyonlu çoklu-iş öğrenme ve çekirdekli bayes çoklu-iş öğrenme, yığıtlı genelleme yöntemiyle birleştiren bir topluluk modeli tasarlanmıştır. Oluşturulan modelin geçerliliğinin onaylanması için ilaç-hücre hattı çiftlerinin aktivite verilerini içeren iki büyük veri kümesi kullanılmıştır. Her iki veri kümesi üzerindeki çapraz doğrulama sonuçları, oluşturulan modelin diğerlerini performans olarak geride bıraktığını göstermektedir. Bu sonuçlar ışığında, orijinal veri setlerinde bulunmayan çiftler için yeni tahminler yapılmıştır. Oluşturulan modelin tahminlerindeki iyileşme, birleştirilecek uygun modelleri seçmekten ve bu modelleri birleştirme yönteminden kaynaklanmaktadır.

Bu yöntem üzerinde çalışılması planlanan bir kaç eklenti daha vardır. İlk olarak, tasarlanan topluluk modeli, farklı öğrenme modelleri seçilerek veya seçilen modelleri birleştirmek için başka yollar kullanılarak genişletilebilir. Belirtildiği gibi, bu seçimler topluluk öğrenme modellerinde kilit rol oynamaktadır.

Bunun yanı sıra, farklı adlandırmalardan kaynaklanan ilaçların veya hücre hatlarının uyuma problemlerini ortadan kaldırdıktan sonra (ör. Hücre hatları eş anlamlıdır ve farklı verisetlerinde farklı olarak adlandırılmıştır.) model, farklı veri setlerini birleştirerek her bir ilaç için daha dengeli veri setleri kullanarak eğitilebilir. Ayrıca, seçilen ve kullanılan öznitelikler, öğrenme modelleri için önemlidir. İlaç tepkisi tahmini için geleneksel yol, gen ifade verilerini kullanılmasıdır ancak bu özellikler ilaçların kimyasal bilgisinin ve diğer genomik özelliklerinin entegrasyonu ile genişletilebilir.

Tez çalışmasında aktarılan diğer yöntem ise ilaçların imzalarının elde edilmesi, imzalardan benzerlik ilişkisinin ortaya konması ve bu benzerlik ilişkilerinden yararlanan bir çoklu-iş öğrenme modelinin oluşturulmasını içerir. İmza benzerlik tabanlı çoklu-iş öğrenme yöntemi olarak adlandırılan bu yöntemin ilaç aktivite tahmini için iki farklı veri kümesiyle yapılan çapraz değerlendirmelerin sonuçlarına bakıldığında, oluşturulan modelin benzerlik ilişkisinin verilmediği modelden daha iyi sonuçlar ortaya koyduğu gözlenmiştir.

LINCS veri kümesinin ve yardımcı bazı tabloların geliştirilmesiyle bu modelin daha da güçleneceği düşünülmektedir. İlaçlar farklı veri kümelerinde farklı kimlikler ile yer aldığı için ilaç imzalarının oluşturulması aşamasında problemlerle karşılaşmıştır. GDSC veri kümesinde bulunan ilaçların yarısından bile azı için imza üretilebilmiştir. İleride bu eşleştirme işlemlerinin daha rahat şekilde yapılabilmesi için gerekli olan arama tablolarının oluşturulmasıyla, daha fazla sayıda ilaç için imza oluşturulacağı düşünülmektedir. Ayrıca LINCS tarafından, veri kümesinde bulunan deneylerin de geliştirilmesiyle bu imzaların güvenilirliği de artacaktır.

GDSC ve CTRP veri kümelerinde yüzlerce hücre hattı bulunurken, LINCS'te üzerinde deney yapılan hücre hattı sayısı şuan 76'dır. İleride diğer hücre hatlarıyla da yapılacak deneylerin eklenmesi daha sağlam bir model oluşturulmasına katkı sağlayacaktır. Bu projenin devamında, ilaç aktivite imzası oluşturma işlemine benzer olarak, hücre aktivite imzasının da oluşturulması ve modelde kullanılması da amaçlanmaktadır. Yapılması planlanan bir diğer çalışma da verilen modellerin bir arada kullanılmasıdır. İmza benzerlik tabanlı regülarizasyonlu çoklu-iş öğrenme modeli, topluluk yöntemi için temel alınan modellerden biri olarak düşünülebilir.

Özetle, bu çalışmada ilaç aktivite tahmini için iki farklı yöntem geliştirilmiştir. İlk çalışma olan topluluk yönteminde, temel alınan modeller yığıtlı genelleme ile birleştirilmiştir. İkinci çalışmada ise ilaç imzaları ile hesaplanan benzerlik ilişkisi kullanılarak regüle edilmiş çoklu-iş öğrenme modeli oluşturulmuştur. Yapılan değerlendirmelerde her iki modelin de ilaç aktivite tahmini için literatüre önemli katkılar yaptığı söylenebilir. Tezin ilk kısmındaki topluluk yönteminin ele alındığı bir makale, uluslararası bir konferansta <sup>3</sup> yayımlanmıştır. Ayrıca her iki çalışma da TÜBİTAK (Proje No: 115E274) tarafından desteklenmektedir.

---

<sup>3</sup>8th International Conference on Knowledge Discovery and Information Retrieval

## KAYNAKLAR

- [1] **Barretina, J., Caponigro, G., Stransky, N., Venkatesan, K., ve diğ.** The Cancer Cell Line Encyclopedia enables predictive modelling of anticancer drug sensitivity. *Nature* 483, 7391 (3 2012), 603–307.
- [2] **Bellmunt, J., Werner, L., Bamias, A., Fay, A. P., ve diğ.** HER2 as a target in invasive urothelial carcinoma. *Cancer Medicine* 4, 6 (6 2015), 844–852.
- [3] **Brem, G. J., Mylonas, I., ve Brüning, A.** Eeyarestatin causes cervical cancer cell sensitization to bortezomib treatment by augmenting ER stress and CHOP expression. *Gynecologic Oncology* 128 (2013), 383–390.
- [4] **Byers, L. A., Diao, L., ve diğ.** An Epithelial-Mesenchymal Transition Gene Signature Predicts Resistance to EGFR and PI3K Inhibitors and Identifies Axl as a Therapeutic Target for Overcoming EGFR Inhibitor Resistance. *Clinical Cancer Research* 19, 1 (1 2013), 279–290.
- [5] **Caruana, R.** Multitask Learning. In *Learning to Learn*. Springer US, Boston, MA, 1998, pp. 95–133.
- [6] **Chan, G. K. Y., Kleinheinz, T. L., ve diğ.** A simple high-content cell cycle assay reveals frequent discrepancies between cell number and ATP and MTS proliferation assays. *PloS one* 8, 5 (2013), e63583.
- [7] **Cortés-Ciriano, I., van Westen, G. J. P., Bouvier, G., ve diğ.** Improved large-scale prediction of growth inhibition patterns using the NCI60 cancer cell line panel. *Bioinformatics* 32, 1 (9 2015), btv529.
- [8] **Costello, J. C., Heiser, L. M., Georgii, E., ve diğ.** A community effort to assess and improve drug sensitivity prediction algorithms. *Nature Biotechnology* 32, 12 (2014), 1202.
- [9] **Cukier, H., Peralta, R., Jin, H., ve diğ.** Preclinical dose scheduling studies of lor-253, a novel anticancer drug, in combination with chemotherapeutics in lung and colon cancers. In *Annual Meeting of the American Association for Cancer Research; Chicago, IL. Philadelphia (PA): AACR; Cancer Res* (2012), vol. 72.
- [10] **DeNardo, D. G., Brennan, D. J., Rexhepaj, E., ve diğ.** Leukocyte Complexity Predicts Breast Cancer Survival and Functionally Regulates Response to Chemotherapy. *Cancer Discovery* 1, 1 (2011).

- [11] **Dietterich, T. G.** Ensemble Methods in Machine Learning. Springer Berlin Heidelberg, 2000, pp. 1–15.
- [12] **Dong, Z., Zhang, N., Li, C., ve diğ.** Anticancer drug sensitivity prediction in cell lines from baseline gene expression through recursive feature selection. *BMC cancer* 15, 1 (2015), 489.
- [13] **Duan, Q., Flynn, C., Niepel, M., Hafner, M., Muhlich, J. L., Fernandez, N. F., Rouillard, A. D., Tan, C. M., Chen, E. Y., Golub, T. R., ve diğ.** Lincs canvas browser: interactive web app to query, browse and interrogate lincs l1000 gene expression signatures. *Nucleic acids research* (2014), gku476.
- [14] **Falgreen, S., Dybkær, K., Young, K. H., Xu-Monette, Z. Y., El-Galaly, T. C., Laursen, M. B., Bødker, J. S., Kjeldsen, M. K., Schmitz, A., Nyegaard, M., Johnsen, H. E., ve Bøgsted, M.** Predicting response to multidrug regimens in cancer patients using cell line experiments and regularised regression models. *BMC Cancer* 15, 1 (12 2015), 235.
- [15] **Fersini, E., Messina, E., Archetti, F., ve diğ.** A p-Median approach for predicting drug response in tumour cells. *BMC Bioinformatics* 15, 1 (12 2014), 353.
- [16] **Friedman, J. H.** Stochastic gradient boosting. *Computational Statistics and Data Analysis* 38, 4 (2002), 367–378.
- [17] **Gholami, A., Hahne, H., Wu, Z., Auer, F., Meng, C., Wilhelm, M., ve Kuster, B.** Global Proteome Analysis of the NCI-60 Cell Line Panel. *Cell Reports* 4, 3 (8 2013), 609–620.
- [18] **Gonen, M., ve Margolin, A. A.** Drug susceptibility prediction against a panel of drugs using kernelized Bayesian multitask learning. *Bioinformatics* 30, 17 (9 2014), i556–i563.
- [19] **Jackson, S. E., ve Chester, J. D.** Personalised cancer medicine. *International Journal of Cancer* 137, 2 (7 2015), 262–266.
- [20] **Jang, I. S., Neto, E. C., Guinney, J., Friend, S. H., ve Margolin, A. a.** Systematic assessment of analytical methods for drug sensitivity prediction from cancer cell line data. *Pacific Symposium on Biocomputing. Pacific Symposium on Biocomputing* (2014), 63–74.
- [21] **Ji, S., ve Ye, J.** An accelerated gradient method for trace norm minimization. In *Proceedings of the 26th annual international conference on machine learning* (2009), ACM, pp. 457–464.
- [22] **Menden, M. P., Iorio, F., Garnett, M., ve diğ.** Machine Learning Prediction of Cancer Cell Sensitivity to Drugs Based on Genomic and Chemical Properties. *PLoS ONE* 8, 4 (4 2013), e61318.
- [23] **National Cancer Database.** American College of Surgeons Commission on Cancer, 2011 Data Submission. *American College of Surgeons* (2013).



- [24] **Neto, E. C., Jang, I. S., Friend, S. H., ve Margolin, A. A.** The Stream algorithm: computationally efficient ridge-regression via Bayesian model averaging, and applications to pharmacogenomic prediction of cancer cell line sensitivity. *Pacific Symposium on Biocomputing. Pacific Symposium on Biocomputing* (2014), 27–38.
- [25] **Nicolau, M., Levine, A. J., ve Carlsson, G.** Topology based data analysis identifies a subgroup of breast cancers with a unique mutational profile and excellent survival. *Proceedings of the National Academy of Sciences* 108, 17 (4 2011), 7265–7270.
- [26] **Qin, Y., Chen, M., ve diğ.** A network flow-based method to predict anticancer drug sensitivity. *PLoS ONE* 10, 5 (2015), 1–14.
- [27] **Rappaport, N., Nativ, N., Stelzer, G., ve diğ.** MalaCards: an integrated compendium for diseases and their annotation. *Database : the journal of biological databases and curation* 2013 (2013), bat018.
- [28] **Rees, M. G., Seashore-Ludlow, B., Cheah, J. H., ve diğ.** Correlating chemical sensitivity and basal gene expression reveals mechanism of action. *Nature Chemical Biology* 12, 2 (12 2015), 109–116.
- [29] **Riddick, G., Song, H., Ahn, S., ve diğ.** Predicting in vitro drug sensitivity using Random Forests. 220–22410.
- [30] **Seashore-Ludlow, B., Rees, M. G., Cheah, J. H., ve diğ.** Harnessing Connectivity in a Large-Scale Small-Molecule Sensitivity Dataset. *Cancer discovery* 5, 11 (11 2015), 1210–23.
- [31] **Sewell, M.** Ensemble learning. *RN* 11, 02.
- [32] **Shi, L., Song, X.-B., Wang, Y., ve diğ.** Docetaxel-conjugated monomethoxy-poly(ethylene glycol)-b-poly(lactide) (mPEG-PLA) polymeric micelles to enhance the therapeutic efficacy in oral squamous cell carcinoma. *RSC Adv.* 6, 49 (2016), 42819–42826.
- [33] **Shoemaker, R. H.** The NCI60 human tumour cell line anticancer drug screen. *Nature Reviews Cancer* 6, 10 (10 2006), 813–823.
- [34] **Siegel, R. L., Miller, K. D., ve Jemal, A.** Cancer statistics, 2016. *CA: A Cancer Journal for Clinicians* 66, 1 (1 2016), 7–30.
- [35] **Tan, M.** Prediction of anti-cancer drug response by kernelized multi-task learning. *Artificial Intelligence in Medicine* 73 (2016), 70–77.
- [36] **Wan, Q., ve Pal, R.** An ensemble based top performing approach for nci-dream drug sensitivity prediction challenge. *PloS one* 9, 6 (2014), e101183.
- [37] **Wolpert, D. H.** Stacked generalization. *Neural Networks* 5, 2 (1992), 241–259.
- [38] **Wu, C., MacLeod, I., ve Su, A. I.** Biogps and mygene. info: organizing online, gene-centric information. *Nucleic acids research* (2012), gks1114.

- [39] **Yang, W., Soares, J., Greninger, P., ve diğ.** Genomics of Drug Sensitivity in Cancer (GDSC): a resource for therapeutic biomarker discovery in cancer cells. *Nucleic Acids Research* 41, D1 (1 2013), D955–D961.
- [40] **Yuan, H., Paskov, I., Paskov, H., González, A. J., ve Leslie, C. S.** Multitask learning improves prediction of cancer drug sensitivity. *Scientific reports* 6 (8 2016), 31619.
- [41] **Zhang, N., Wang, H., Fang, Y., ve diğ.** Predicting Anticancer Drug Responses Using a Dual-Layer Integrated Cell Line-Drug Network Model. *PLoS computational biology* 11, 9 (2015), e1004498.
- [42] **Zhao, Z., Fu, G., Liu, S., ve diğ.** Drug activity prediction using multiple-instance learning via joint instance and feature selection. *BMC Bioinformatics* 14 Suppl 1, Suppl 14 (2013), S16.
- [43] **Zhou, J.** Multi-task learning in crisis event classification. Tech. rep., Tech. Rep., [http://www. public. asu. edu/jzhou29](http://www.public.asu.edu/jzhou29).
- [44] **Zhou, J., Chen, J., ve Ye, J.** User’s Manual MALSAR: Multi-tAsk Learning via StructurAl Regularization. *Arizona State University* (2012).

## **EKLER**

**EK 1** : Şekil 5.1 ve 5.2 için İlaç Listesi

**EK 2** : Şekil 5.3 için İlaç Listesi

**EK 3** : Şekil 5.4 için İlaç Listesi

**EK 4** : Şekil 5.5 için İlaç Listesi

## EK 1

1. Erlotinib	54. AZD6482	107. UNC0638	161. QL-X-138	215. PD-0325901
2. Rapamycin	55. JNK-9L	108. XL-184	162. XMD15-27	216. SB590885
3. Sunitinib	56. PF-562271	109. WZ3105	163. T0901317	217. selumetinib
4. PHA-665752	57. HG-6-64-1	110. XMD14-99	164. EX-527	218. AZD6482
5. MG-132	58. JQ1	111. AC220	165. THZ-2-49	219. CCT007093
6. Paclitaxel	59. JQ12	112. CP724714	166. KIN001-270	220. EHT 1864
7. Cyclopamine	60. DMOG	113. JW-7-24-1	167. THZ-2-102-1	221. BMS-708163
8. AZ628	61. FTI-277	114. NPK76-II-72-1	168. AICAR	222. BMS-536924
9. Sorafenib	62. OSU-03012	115. STF-62247	169. Camptothecin	223. Cetuximab
10. VX-680	63. Shikonin	116. NG-25	170. Vinblastine	224. PF-4708671
11. Imatinib	64. AKT inhibitor VIII	117. TL-1-85	171. Cisplatin	225. JNJ-26854165
12. TAE684	65. Embelin	118. VX-11e	172. Cytarabine	226. HG-5-113-01
13. Crizotinib	66. FH535	119. FR-180204	173. Docetaxel	227. HG-5-88-01
14. Saracatinib	67. PAC-1	120. Tubastatin A	174. Methotrexate	228. TW 37
15. S-Trityl-L-cysteine	68. IPA-3	121. Zibotentan	175. ATRA	229. XMD111-85h
16. Z-LLNle-CHO	69. GSK-650394	122. YM155	176. Gefitinib	230. ZG-10
17. Dasatinib	70. BAY 61-3606	123. NSC-207895	177. Navitoclax	231. XMD8-92
18. GNF-2	71. 5-Fluorouracil	124. VNLG/124	178. Vorinostat	232. QL-VIII-58
19. CGP-60474	72. Thapsigargin	125. AR-42	179. Nilotinib	233. CCT018159
20. CGP-082996	73. Obatoclox Mesylate	126. CUDC-101	180. RDEA119	234. AG-014699
21. A-770041	74. BMS-754807	127. Belinostat	181. CI-1040	235. GSK269962A
22. WH-4-023	75. Lisitinib	128. I-BET-762	182. Temsirolimus	236. SB 505124
23. WZ-1-84	76. Bexarotene	129. CAY10603	183. Olaparib	237. Tamoxifen
24. BI-2536	77. Bleomycin	130. Linifanib	184. Veliparib	238. QL-XII-61
25. BMS-536924	78. LFM-A13	131. BIX02189	185. Bosutinib	239. JQ1
26. BMS-509744	79. GW-2580	132. CH5424802	186. Lenalidomide	240. PFI-1
27. CMK	80. AUY922	133. EKB-569	187. Axitinib	241. IOX2
28. Pyrimethamine	81. Phenformin	134. GSK2126458	188. AZD7762	242. UNC0638
29. JW-7-52-1	82. Bryostatins 1	135. KIN001-236	189. GW 441756	243. YK 4-279
30. A-443654	83. Pazopanib	136. KIN001-244	190. CEP-701	244. CHIR-99021
31. GW843682X	84. LAQ824	137. KIN001-055	191. SB 216763	245. (5Z)-7-Oxozeanol
32. MS-275	85. Epothilone B	138. KIN001-260	192. 17-AAG	246. piperlongumine
33. Parthenolide	86. GSK1904529A	139. KIN001-266	193. VX-702	247. FK866
34. KIN001-135	87. BMS345541	140. Masitinib	194. AMG-706	248. Talazoparib
35. TGX221	88. Tipifarnib	141. MP470	195. KU-55933	249. rTRAIL
36. Bortezomib	89. BMS-708163	142. MPS-1-IN-1	196. Elesclomol	250. UNC1215
37. XMD8-85	90. Ruxolitinib	143. BHG712	197. Afatinib	251. SGC0946
38. Roscovitine	91. AS601245	144. OSI-930	198. GDC0449	252. XAV939
39. Salubrinal	92. Ispinesib Mesylate	145. OSI-027	199. PLX4720	253. PLX4720
40. Lapatinib	93. TL-2-105	146. CX-5461	200. BX-795	254. Trametinib
41. GSK269962A	94. AT-7519	147. PHA-793887	201. NU-7441	255. Dabrafenib
42. Doxorubicin	95. TAK-715	148. PI-103	202. SL 0101-1	256. Temozolomide
43. Etoposide	96. BX-912	149. PIK-93	203. BIRB 0796	257. Afatinib
44. Gemcitabine	97. ZSTK474	150. SB52334	204. JNK Inhibitor VIII	258. Bleomycin (50 uM)
45. Mitomycin C	98. AS605240	151. TPCA-1	205. 681640	259. SN-38
46. Vinorelbine	99. Genentech Cpd 10	152. TG101348	206. Nutlin-3a (-)	260. Olaparib
47. NSC-87877	100. GSK1070916	153. Foretinib	207. PD-173074	261. selumetinib
48. Bicalutamide	101. KIN001-102	154. Y-39983	208. ZM-447439	262. Bicalutamide
49. QS11	102. LY317615	155. YM201636	209. RO-3306	263. RDEA119
50. CP466722	103. GSK429286A	156. Tivozanib	210. MK-2206	264. GDC0941
51. Midostaurin	104. FMK	157. GSK690693	211. PD-0332991	265. MLN4924
52. CHIR-99021	105. QL-XII-47	158. SNX-2112	212. BEZ235	
53. AP-24534	106. CAL-101	159. QL-XI-92	213. GDC0941	
		160. XMD13-2	214. AZD8055	

## EK 2

1. BRD-K46556387	58. BRD-K52075040	115. BRD-K90382497	172. BRD-K11172604	229. BRD-K89162000
2. BRD-K86574132	59. BRD-K73397362	116. BRD-K77625799	173. BRD-K81651477	230. BRD-K55478147
3. BRD-K35716340	60. BRD-K04923131	117. BRD-K23984367	174. BRD-K61829047	231. BRD-K17349619
4. BRD-K89692698	61. BRD-K49328571	118. BRD-K59962020	175. BRD-K12502280	232. BRD-K39706510
5. BRD-A18763547	62. BRD-K64052750	119. BRD-A94153989	176. BRD-K67844266	233. BRD-K45748132
6. BRD-K89329876	63. BRD-K70401845	120. BRD-K90370028	177. BRD-K52836380	234. BRD-K49456190
7. BRD-K64634304	64. BRD-K94991378	121. BRD-K71281111	178. BRD-K32330832	235. BRD-A68631409
8. BRD-K45401373	65. BRD-K26818574	122. BRD-K82746043	179. BRD-K80672993	236. BRD-K86856088
9. BRD-K15563106	66. BRD-K52313696	123. BRD-K92441787	180. BRD-K53855319	237. BRD-A28105619
10. BRD-A93255169	67. BRD-K02130563	124. BRD-K81528515	181. BRD-K61662457	238. BRD-K50799972
11. BRD-K19295594	68. BRD-K75430629	125. BRD-K42828737	182. BRD-K30064966	239. BRD-K17060750
12. BRD-K29458283	69. BRD-A38030642	126. BRD-K17068645	183. BRD-K89391146	240. BRD-K63923597
13. BRD-K24844714	70. BRD-K81418486	127. BRD-K76674262	184. BRD-K92991072	241. BRD-K99749624
14. BRD-K03406345	71. BRD-K84937637	128. BRD-K28907958	185. BRD-K98538768	242. BRD-K87142802
15. BRD-K89732114	72. BRD-K25311561	129. BRD-K17140735	186. BRD-K47335880	243. BRD-K19540840
16. BRD-A28746609	73. BRD-K06854232	130. BRD-K06593056	187. BRD-K10466330	244. BRD-K66175015
17. BRD-K93754473	74. BRD-K05653692	131. BRD-K20755323	188. BRD-K34581968	245. BRD-K86930074
18. BRD-A35588707	75. BRD-K71935468	132. BRD-K88544581	189. BRD-K15108141	246. BRD-K50168500
19. BRD-K50128260	76. BRD-K60230970	133. BRD-A36318220	190. BRD-K19894101	247. BRD-K15600710
20. BRD-K22134346	77. BRD-K77908580	134. BRD-K10705233	191. BRD-K02113016	248. BRD-K71035033
21. BRD-K02407574	78. BRD-K17743125	135. BRD-K93176058	192. BRD-A81541225	249. BRD-K63068307
22. BRD-K13032584	79. BRD-K64606589	136. BRD-K00317371	193. BRD-K17610631	250. BRD-K28428262
23. BRD-K74148702	80. BRD-K04466929	137. BRD-K37865504	194. BRD-K02492147	251. BRD-K67578145
24. BRD-K13044802	81. BRD-A94377914	138. BRD-K47150025	195. BRD-K25987073	252. BRD-K50140147
25. BRD-K55591206	82. BRD-A42556028	139. BRD-K25737009	196. BRD-K93095519	253. BRD-K43389698
26. BRD-K43149758	83. BRD-K82109576	140. BRD-K51575138	197. BRD-K09587429	254. BRD-K16478699
27. BRD-K59456551	84. BRD-K33106058	141. BRD-K37392901	198. BRD-K44741158	255. BRD-A41692738
28. BRD-A70155556	85. BRD-A34817987	142. BRD-K96799727	199. BRD-K00088062	256. BRD-K67566344
29. BRD-K12994359	86. BRD-A25004090	143. BRD-A36630025	200. BRD-K62965247	257. BRD-K64785675
30. BRD-A09722536	87. BRD-K38985961	144. BRD-K80183349	201. BRD-K63195589	258. BRD-K06750613
31. BRD-K35520305	88. BRD-K11740178	145. BRD-K66532283	202. BRD-K67112618	259. BRD-K02965346
32. BRD-K35960502	89. BRD-A91658086	146. BRD-K68065987	203. BRD-K38264551	260. BRD-K88742110
33. BRD-K19352500	90. BRD-K92856060	147. BRD-K05402890	204. BRD-A22997170	261. BRD-K01436366
34. BRD-A67097164	91. BRD-K10882151	148. BRD-K66453893	205. BRD-K20285085	262. BRD-K22828899
35. BRD-K92093830	92. BRD-K45681478	149. BRD-K11533227	206. BRD-K76703230	263. BRD-K76964878
36. BRD-K35708212	93. BRD-A83255679	150. BRD-K14844214	207. BRD-K13049116	264. BRD-K44847641
37. BRD-K05649647	94. BRD-A36275421	151. BRD-K62801835	208. BRD-K04623885	265. BRD-K36852164
38. BRD-K24132293	95. BRD-K05392795	152. BRD-K41597374	209. BRD-K92241597	266. BRD-K23853216
39. BRD-K26531177	96. BRD-K29711668	153. BRD-K63431240	210. BRD-K85606544	267. BRD-K69982010
40. BRD-A80213327	97. BRD-K20215950	154. BRD-K13999467	211. BRD-K73261812	268. BRD-K24376488
41. BRD-K01567962	98. BRD-K64497429	155. BRD-A31107743	212. BRD-K68548958	269. BRD-K36739687
42. BRD-K08608874	99. BRD-K59204667	156. BRD-K32536677	213. BRD-K60219430	270. BRD-K35498412
43. BRD-K55696337	100. BRD-K61128465	157. BRD-K87317732	214. BRD-K58550667	271. BRD-K01877528
44. BRD-K71726959	101. BRD-K44241590	158. BRD-K83213911	215. BRD-K54256913	272. BRD-K35723520
45. BRD-A02481876	102. BRD-A86956638	159. BRD-A13122391	216. BRD-K53414658	273. BRD-A56592690
46. BRD-K37798499	103. BRD-K45841694	160. BRD-A12230535	217. BRD-K64642496	274. BRD-K09485525
47. BRD-K15318909	104. BRD-A35108200	161. BRD-K07303502	218. BRD-K23547378	275. BRD-K83336168
48. BRD-K81473043	105. BRD-A40802033	162. BRD-K98404142	219. BRD-K71781559	276. BRD-K28360340
49. BRD-K26664453	106. BRD-K92723993	163. BRD-K09907482	220. BRD-K04800985	277. BRD-K69932463
50. BRD-K68143200	107. BRD-K92041145	164. BRD-K08417745	221. BRD-K19103580	278. BRD-K94145482
51. BRD-K78599730	108. BRD-K79254416	165. BRD-K53903639	222. BRD-K64610608	279. BRD-K01121114
52. BRD-K59670716	109. BRD-K16189898	166. BRD-K40853697	223. BRD-K02251932	280. BRD-A67788537
53. BRD-K69608737	110. BRD-K64890080	167. BRD-K88510285	224. BRD-K55116708	281. BRD-A63646118
54. BRD-K66874953	111. BRD-K61166597	168. BRD-K56301217	225. BRD-K16147474	282. BRD-K54606188
55. BRD-K67298865	112. BRD-K29905972	169. BRD-A01145011	226. BRD-K33199242	283. BRD-K51490254
56. BRD-K61323504	113. BRD-K57080016	170. BRD-K67868012	227. BRD-K52037352	284. BRD-K85133207
57. BRD-K26603252	114. BRD-K12184916	171. BRD-K00627859	228. BRD-K37390332	285. BRD-K13566078

291. BRD-K16730910	321. BRD-K48477130	351. BRD-K99655327	381. BRD-K74065929	411. BRD-K49215523
292. BRD-K81728688	322. BRD-K92428232	352. BRD-K45935533	382. BRD-K21718444	412. BRD-K75664313
293. BRD-K25340465	323. BRD-K28456706	353. BRD-K35604418	383. BRD-K93123848	413. BRD-K64117221
294. BRD-K56277358	324. BRD-K22477529	354. BRD-K62358710	384. BRD-K81458380	414. BRD-K50501969
295. BRD-K39974922	325. BRD-K93367411	355. BRD-K08109215	385. BRD-K87737963	415. BRD-K42137908
296. BRD-A62182663	326. BRD-K84964099	356. BRD-K92960067	386. BRD-K60866521	416. BRD-K99006945
297. BRD-K83289131	327. BRD-K96335988	357. BRD-K15179513	387. BRD-K43644456	417. BRD-K33379087
298. BRD-K09499853	328. BRD-A02303741	358. BRD-K34022604	388. BRD-K87124298	418. BRD-K12343256
299. BRD-K06792661	329. BRD-K81491172	359. BRD-K75308783	389. BRD-K02017404	419. BRD-K22024824
300. BRD-K78659596	330. BRD-A34462049	360. BRD-K65928735	390. BRD-K29395450	420. BRD-K13662825
301. BRD-K69840642	331. BRD-K79877282	361. BRD-A82720512	391. BRD-K86525559	421. BRD-K43410529
302. BRD-K19687926	332. BRD-K79239947	362. BRD-K47981327	392. BRD-K18589165	422. BRD-K51831558
303. BRD-K29313308	333. BRD-K31514534	363. BRD-K59437938	393. BRD-K29968218	423. BRD-K87774949
304. BRD-K74236984	334. BRD-K78667050	364. BRD-A04287157	394. BRD-A09890259	424. BRD-K90860366
305. BRD-K19477839	335. BRD-K81335284	365. BRD-K04905989	395. BRD-K78978711	425. BRD-K42260513
306. BRD-K19796430	336. BRD-K03391209	366. BRD-A52530402	396. BRD-K54640016	426. BRD-K68682020
307. BRD-K75295174	337. BRD-A76527075	367. BRD-K28392481	397. BRD-K51544265	427. BRD-K84754008
308. BRD-K58772419	338. BRD-K66430217	368. BRD-A30032852	398. BRD-M00053801	428. BRD-A86708339
309. BRD-K53972329	339. BRD-K27624156	369. BRD-A12571627	399. BRD-K71467466	429. BRD-K37764012
310. BRD-A42083487	340. BRD-K59146805	370. BRD-K18850819	400. BRD-K90570971	430. BRD-K09951645
311. BRD-K52911425	341. BRD-K31406481	371. BRD-K65592642	401. BRD-K14870255	431. BRD-K55125219
312. BRD-K37720887	342. BRD-K97651142	372. BRD-A05715709	402. BRD-K70301465	432. BRD-K58435339
313. BRD-K64800655	343. BRD-K93442924	373. BRD-K90746395	403. BRD-K16761703	433. BRD-A68258842
314. BRD-K93918653	344. BRD-K71487808	374. BRD-K08177763	404. BRD-K76894955	434. BRD-K83607951
315. BRD-K88025533	345. BRD-A71883111	375. BRD-A05821830	405. BRD-A59431241	435. BRD-K34222889
316. BRD-K74514084	346. BRD-K50145167	376. BRD-K84466663	406. BRD-K62391742	436. BRD-K43428468
317. BRD-K72420232	347. BRD-K60750172	377. BRD-K55187425	407. BRD-K24556407	437. BRD-A15100685
318. BRD-K24690302	348. BRD-K36363294	378. BRD-K12040459	408. BRD-K57261999	438. BRD-K08589866
319. BRD-K62825658	349. BRD-K27955832	379. BRD-A57798112	409. BRD-K52560704	439. BRD-K02834582
320. BRD-K56343971	350. BRD-K11593101	380. BRD-K36016295	410. BRD-K54997624	

## EK 3

1. A-443654	28. Cisplatin	55. JNK-9L	82. PD-0325901	109. TAK-715
2. AMG-706	29. Crizotinib	56. JW-7-24-1	83. PD-0332991	110. TGX221
3. AP-24534	30. Cyclopamine	57. JW-7-52-1	84. PD-173074	111. TPCA-1
4. AT-7519	31. Cytarabine	58. KIN001-055	85. PF-562271	112. TW 37
5. ATRA	32. Dabrafenib	59. KIN001-135	86. PHA-665752	113. Tamoxifen
6. AUY922	33. Dasatinib	60. KIN001-244	87. PHA-793887	114. Temozolomide
7. AZ628	34. Docetaxel	61. KIN001-266	88. PI-103	115. Temsirolimus
8. AZD6482	35. Doxorubicin	62. KU-55933	89. PLX4720	116. Thapsigargin
9. Afatinib	36. EX-527	63. LAQ824	90. Paclitaxel	117. Tipifarnib
10. Axitinib	37. Elesclomol	64. LFM-A13	91. Parthenolide	118. Tivozanib
11. BI-2536	38. Embelin	65. Lapatinib	92. Pazopanib	119. Trametinib
12. BIRB0796	39. Epothilone B	66. Lenalidomide	93. Phenformin	120. VX-680
13. BMS-509744	40. Erlotinib	67. Liniifanib	94. Pyrimethamine	121. VX-702
14. BMS-536924	41. Etoposide	68. MG-132	95. QL-X-138	122. Veliparib
15. BMS-754807	42. FH535	69. MK-2206	96. QL-XI-92	123. Vinblastine
16. BX-795	43. FR-180204	70. MS-275	97. QL-XII-47	124. Vinorelbine
17. BX-912	44. Foretinib	71. Masitinib	98. RO-3306	125. Vorinostat
18. Belinostat	45. GDC0941	72. Methotrexate	99. Rapamycin	126. WH-4-023
19. Bicalutamide	46. GNF-2	73. Midostaurin	100. Roscovitine	127. XMD8-85
20. Bortezomib	47. GSK-650394	74. NPK76-II-72-1	101. Ruxolitinib	128. Y-39983
21. Bosutinib	48. GW843682X	75. NU-7441	102. SB590885	129. ZG-10
22. CEP-701	49. Gefitinib	76. Navitoclax	103. SN-38	130. ZM-447439
23. CGP-60474	50. Gemcitabine	77. Nilotinib	104. STF-62247	131. Zibotentan
24. CHIR-99021	51. HG-5-113-01	78. OSI-027	105. Salubrinal	132. piperlongumine
25. CI-1040	52. HG-5-88-01	79. OSI-930	106. Saracatinib	133. selumetinib
26. CP466722	53. Imatinib	80. Olaparib	107. Sorafenib	
27. Camptothecin	54. JNJ-26854165	81. PAC-1	108. Sunitinib	

## EK 4

1. BRD-A01145011	59. BRD-K06593056	117. BRD-K27224038	175. BRD-K55187425	233. BRD-K75308783
2. BRD-A02481876	60. BRD-K06750613	118. BRD-K27986637	176. BRD-K55478147	234. BRD-K75430629
3. BRD-A09722536	61. BRD-K06792661	119. BRD-K28360340	177. BRD-K55591206	235. BRD-K76674262
4. BRD-A12230535	62. BRD-K06854232	120. BRD-K28428262	178. BRD-K55696337	236. BRD-K76703230
5. BRD-A13122391	63. BRD-K07303502	121. BRD-K28907958	179. BRD-K56277358	237. BRD-K76964878
6. BRD-A15100685	64. BRD-K08417745	122. BRD-K29313308	180. BRD-K56343971	238. BRD-K77908580
7. BRD-A18763547	65. BRD-K08589866	123. BRD-K29458283	181. BRD-K56343971	239. BRD-K78431006
8. BRD-A22997170	66. BRD-K09485525	124. BRD-K29905972	182. BRD-K57080016	240. BRD-K78599730
9. BRD-A25004090	67. BRD-K09499853	125. BRD-K30064966	183. BRD-K58550667	241. BRD-K78659596
10. BRD-A28105619	68. BRD-K09778810	126. BRD-K32330832	184. BRD-K58772419	242. BRD-K79254416
11. BRD-A28746609	69. BRD-K09907482	127. BRD-K32536677	185. BRD-K59369769	243. BRD-K79877282
12. BRD-A31107743	70. BRD-K09951645	128. BRD-K33106058	186. BRD-K59456551	244. BRD-K80672993
13. BRD-A34462049	71. BRD-K10466330	129. BRD-K33379087	187. BRD-K59962020	245. BRD-K81418486
14. BRD-A34817987	72. BRD-K10705233	130. BRD-K33514849	188. BRD-K60219430	246. BRD-K81473043
15. BRD-A35108200	73. BRD-K10882151	131. BRD-K33583600	189. BRD-K60230970	247. BRD-K81528515
16. BRD-A35588707	74. BRD-K11533227	132. BRD-K34157611	190. BRD-K60866521	248. BRD-K81651477
17. BRD-A36275421	75. BRD-K12040459	133. BRD-K34581968	191. BRD-K61323504	249. BRD-K82109576
18. BRD-A36318220	76. BRD-K12184916	134. BRD-K35520305	192. BRD-K61662457	250. BRD-K82746043
19. BRD-A36630025	77. BRD-K12343256	135. BRD-K35708212	193. BRD-K61829047	251. BRD-K83213911
20. BRD-A38030642	78. BRD-K12502280	136. BRD-K35716340	194. BRD-K62965247	252. BRD-K83289131
21. BRD-A41692738	79. BRD-K12994359	137. BRD-K35723520	195. BRD-K63068307	253. BRD-K83336168
22. BRD-A56592690	80. BRD-K13032584	138. BRD-K35960502	196. BRD-K63195589	254. BRD-K84937637
23. BRD-A59431241	81. BRD-K13044802	139. BRD-K37392901	197. BRD-K63431240	255. BRD-K85133207
24. BRD-A62182663	82. BRD-K13049116	140. BRD-K37720887	198. BRD-K63923597	256. BRD-K85606544
25. BRD-A63646118	83. BRD-K13566078	141. BRD-K37798499	199. BRD-K64052750	257. BRD-K86535717
26. BRD-A67097164	84. BRD-K13999467	142. BRD-K37865504	200. BRD-K64606589	258. BRD-K86574132
27. BRD-A67788537	85. BRD-K14844214	143. BRD-K38985961	201. BRD-K64610608	259. BRD-K86930074
28. BRD-A68631409	86. BRD-K15108141	144. BRD-K40853697	202. BRD-K64634304	260. BRD-K87142802
29. BRD-A70155556	87. BRD-K15318909	145. BRD-K43149758	203. BRD-K64642496	261. BRD-K87737963
30. BRD-A75817871	88. BRD-K15563106	146. BRD-K43389698	204. BRD-K64785675	262. BRD-K88510285
31. BRD-A76279427	89. BRD-K15600710	147. BRD-K44241590	205. BRD-K64800655	263. BRD-K88544581
32. BRD-A80213327	90. BRD-K16189898	148. BRD-K45401373	206. BRD-K64890080	264. BRD-K88742110
33. BRD-A81541225	91. BRD-K16478699	149. BRD-K45681478	207. BRD-K66175015	265. BRD-K89162000
34. BRD-A83255679	92. BRD-K17140735	150. BRD-K47150025	208. BRD-K66874953	266. BRD-K89329876
35. BRD-A93255169	93. BRD-K17349619	151. BRD-K47335880	209. BRD-K67298865	267. BRD-K89391146
36. BRD-A94377914	94. BRD-K17743125	152. BRD-K47503321	210. BRD-K67566344	268. BRD-K89692698
37. BRD-K00088062	95. BRD-K17953061	153. BRD-K49075727	211. BRD-K67578145	269. BRD-K89732114
38. BRD-K00317371	96. BRD-K19295594	154. BRD-K49290616	212. BRD-K67844266	270. BRD-K90370028
39. BRD-K00627859	97. BRD-K19352500	155. BRD-K49328571	213. BRD-K67868012	271. BRD-K90382497
40. BRD-K01121114	98. BRD-K19416115	156. BRD-K49456190	214. BRD-K68065987	272. BRD-K92093830
41. BRD-K01436366	99. BRD-K19477839	157. BRD-K50128260	215. BRD-K68143200	273. BRD-K92241597
42. BRD-K01567962	100. BRD-K19540840	158. BRD-K50140147	216. BRD-K68548958	274. BRD-K92428232
43. BRD-K01737880	101. BRD-K19687926	159. BRD-K50168500	217. BRD-K69608737	275. BRD-K92723993
44. BRD-K01877528	102. BRD-K19796430	160. BRD-K50799972	218. BRD-K69840642	276. BRD-K92991072
45. BRD-K02113016	103. BRD-K19894101	161. BRD-K51490254	219. BRD-K69932463	277. BRD-K93095519
46. BRD-K02130563	104. BRD-K20285085	162. BRD-K51575138	220. BRD-K70401845	278. BRD-K93176058
47. BRD-K02407574	105. BRD-K20755323	163. BRD-K51781482	221. BRD-K71035033	279. BRD-K93367411
48. BRD-K02492147	106. BRD-K22134346	164. BRD-K52075040	222. BRD-K71726959	280. BRD-K93754473
49. BRD-K02526760	107. BRD-K22828899	165. BRD-K52313696	223. BRD-K71935468	281. BRD-K93918653
50. BRD-K02965346	108. BRD-K23984367	166. BRD-K52836380	224. BRD-K72264770	282. BRD-K94991378
51. BRD-K03406345	109. BRD-K24132293	167. BRD-K52911425	225. BRD-K72420232	283. BRD-K96431673
52. BRD-K03449891	110. BRD-K24690302	168. BRD-K53414658	226. BRD-K73261812	284. BRD-K96799727
53. BRD-K04466929	111. BRD-K24844714	169. BRD-K53792571	227. BRD-K73397362	285. BRD-K98404142
54. BRD-K04623885	112. BRD-K25311561	170. BRD-K53855319	228. BRD-K74065929	286. BRD-K99749624
55. BRD-K04923131	113. BRD-K25340465	171. BRD-K53903639	229. BRD-K74148702	
56. BRD-K05402890	114. BRD-K25737009	172. BRD-K53972329	230. BRD-K74236984	
57. BRD-K05649647	115. BRD-K26664453	173. BRD-K54256913	231. BRD-K74514084	
58. BRD-K05653692	116. BRD-K26818574	174. BRD-K55116708	232. BRD-K75295174	



## ÖZGEÇMİŞ

**Ad-Soyad** : Ertan Tolan  
**Uyruđu** : T.C.  
**Dođum Tarihi ve Yeri** : 18.03.1991 - Aydın  
**E-posta** : e.tolan@etu.edu.tr

### ÖĐRENİM DURUMU:

- **Lisans** : 2014, TOBB ETU, Bilgisayar Mühendisliđi

### MESLEKİ DENEYİM VE ÖDÜLLER:

Yıl	Yer	Görev
2014-2016	TOBB ETU	Burslu Yüksek Lisans Öğrencisi

### YABANCI DİL: İNGİLİZCE

### TEZDEN TÜRETİLEN YAYINLAR, SUNUMLAR VE PATENTLER:

- E. Tolan and M. Tan. Anti-cancer drug activity prediction by ensemble learning. In International Conference on Knowledge Discovery and Information Retrieval, 2016.