



Erciyes University Journal of the Institute of Science and Technology
Erciyes Üniversitesi Fen Bilimleri Enstitüsü Dergisi
 ISSN 1012-2354



Cilt (Volume): 28, Sayı (Issue): 1, Ocak/January-2012
<http://fbe.erciyes.edu.tr/>

HIV-1 Proteaz Özgünlüğünün Yeni Bir Öznitelik Temsili Yöntemi ile Proteomik Analizi

*Murat Gök¹, Ahmet Turan Özcerit²

¹Department of Computer Engineering, Yalova University, Turkey

²Department of Computer & Electronics, Sakarya University, Sakarya Turkey

ÖZET

Bu çalışmada HIV-1 proteaz enzimi bölünme kısımlarının tahmini için Fizikokimyasal Tabanlı Kodlama Yöntemi (FTKY) adı verilen yeni bir öznitelik kodlama yöntemi uygulandı. FTKY, seçilen en iyi 10-fk (fizikokimyasal), 20-fk, 30-fk, 40-fk ve 50-fk özelliğe göre sınıf doğruluğu, duyarlık ve Alıcı İşletim Karakteristiği Eğrisi Altında Kalan Alan (AİKAA) değerleri bakımından Doğrusal Destek Vektör Makineleri (DDVM) yöntemi kullanılarak test edilmiştir. Testlerde güncel iki HIV-1 proteaz veri seti, PR-1625 ve PR-3261 kullanılmıştır. Elde edilen deneysel sonuçlara göre 10-fk'ya göre yapılan kodlamalarda, PR-1625 veri seti üzerinde en yüksek performans elde edilirken PR-3261'de ise en düşük performans elde edilmiştir. Elde edilen deneysel sonuçlara göre FTKY, tek sınıflandırıcı üzerinde PR-1625 üzerinde en yüksek sınıf doğruluğunu % 95,21 ile en iyi 10-fk, PR-3261 üzerinde ise % 94,37 sınıf doğruluğu ile en iyi 30-fk vermiştir.

Anahtar Kelimeler:

HIV-1 Proteaz
 Özgünlüğü, Öznitelik
 Temsili,
 Peptit Sınıflandırma

Proteomic Analysis of HIV-1 Protease Specificity With A New Feature Encoding Method

ABSTRACT

In this study a new feature encoding scheme named FTKY (Physicochemical Based Encoding Method) has been developed for HIV-1 protease site prediction. FTKY has been tested according to selected best 10-pc(physicochemical), 20-pc, 30-pc, 40-pc and 50-pc by means of accuracy, specificity and AUROC (Area Under Receiver Operating Characteristic Curve) on Linear Support Vector Machines. Tests have been conducted on two up-to-date HIV-1 protease datasets, PR-1625 and PR-3265. According to empirical results, FTKY has been performed better prediction of accuracy 95.21 % on PR-1625 according to best 10-pc and accuracy 94.37 % when using PR-3261 according to best 30-pc on a standalone classifier.

Keywords:
 HIV-1 Protease
 Specificity, Feature
 Representation,
 Peptide
 Classification

1. Giriş

HIV-1 proteaz, AIDS virüsünün çoğalması için hayati önemi olan bir enzimdir [1]. Virüsün çoğalması ancak HIV-1 proteazın çoklu proteinleri uygun konumlarından kesmesi ile mümkündür [2]. HIV-1 proteaz, proteinlerin kesme konumlarını belirleme işlemini bir mekanizma içinde yapmaktadır. Eğer bu mekanizmanın şifreleri çözümlerse, uygun baskılayıcı ilaç geliştirilebilir. Böylece baskılayıcı ilaç molekülleri proteazın aktif (katalitik) bölgesine sızabilirler ve aktif bölgeyi tıkayarak proteaz enziminin kesme fonksiyonunu yerine getirmesine mani olurlar [3]. Enzim – baskılayıcı ilaç ilişkisini, anahtar – kilit örneğine benzetebiliriz. Uygun anahtarın kilidi devreye soktuğu gibi uygun baskılayıcı ilaçlar da enzimin işlevini yerine getirmesine mani olurlar. Bu nedenle proteaz enziminin kesilme konumlarının doğru tahmini baskılayıcı ilaç geliştirilmesi ve dolayısıyla AIDS hastalığının durdurulması açısından hayati önem taşımaktadır.

Geçmiş yıllarda HIV-1 protease kesme konumlarının tahmininde makine öğrenmesi temelli çalışmalar yapılmıştır. Bu problemi çözmek için [4]'de yapay sinir ağı yöntemi uygulanmıştır. [5]'de HIV-1 özgünlüğü probleminin doğrusal bir problem olduğu ve DDVM yönteminin, problemin çözümü için en doğru yöntem olduğu sonucuna varılmıştır. [6]'de Quasi kalıntı yöntemi, sınıflandırıcılar birleştirilerek (fusion) uygulanmış ve % 97.5 sınıfı doğruluğu elde edilmiştir. Fakat yapılan bu çalışmada veri manipülasyonu yapıldığı bildirilmektedir [7]. Ayrıca tüm bu çalışmalarda 362 amino asitten oluşan, güncel olmayan bir veri seti kullanılmıştır [4]. Bu veri seti, içerdiği peptit sayısı açısından sınıflandırma algoritmaları için yetersizdir ve günümüzde geçerliliğini yitirmiştir.

Bu çalışmada, literatürde yer alan amino asitlerin 544 fizikokimyasal özellikleri [8] göz önünde bulundurularak FTKY adı verilen bir öznitelik kodlama yöntemi geliştirilmiş ve güncel iki HIV-1 proteaz veri seti kullanılarak HIV-1 proteazın proteinleri kesme konumlarını tahmini problemine uygulanmıştır.

2. Gereç ve Yöntemler

2.1. Peptit Dizilimi

Bir protein doğada bulunan 20 amino asidin, $\mathcal{A} = \{s_1, s_2, \dots, s_{20}\} = \{A, C, D, E, F, G, H, I, K, L, M, N, P, Q, R, S, T, V, W, Y\}$ çeşitli kombinasyonlarda dizilimi ile meydana gelir. Protein içindeki aktif kısımlar, $P = P_4 P_3 P_2 P_1 \downarrow P_1 P_2 P_3 P_4$ ile ifade edilen sekizli amino asit peptit dizilimlerinden oluşur. Bu ifade de \downarrow simgesi P_1 ile P_1 arasında bir makas bağ olduğunu belirtmektedir ve peptit bu kısımdan kesilmektedir [9]. HIV-1 enzimi bölünme kısımlarının

tahmini problemi, sekiz amino asitten oluşan P dizilimlerinin kesilmiş (cleavage) sınıfa mı, kesilmemiş (noncleavage) sınıfa mı ait olduğunun tespit edilmeye çalışıldığı sayısal bir sınıflandırma problemidir.

2.2. PR-1625 ve PR-3261 HIV-1 Proteaz Veri Setleri

HIV-1 proteaz/substrat etkileşimi için Kontijevskis [10] tarafından 2007 yılında 1625 peptit diziliminden oluşan bir veri seti (PR-1625) yayınlanmıştır. Bu örüntü verilerinin 374'ü kesilmiş, 1251'i kesilmemiş peptittir. 2008 yılında Schilling [11] tarafından daha geniş bir veri seti (PR-3261) yayınlanmıştır. PR-3261 veri seti, 436 kesilmiş, 2825 kesilmemiş peptitten oluşmaktadır. Her iki veri seti arasında % 7 oranında küçük bir benzerlik bulunmaktadır [12].

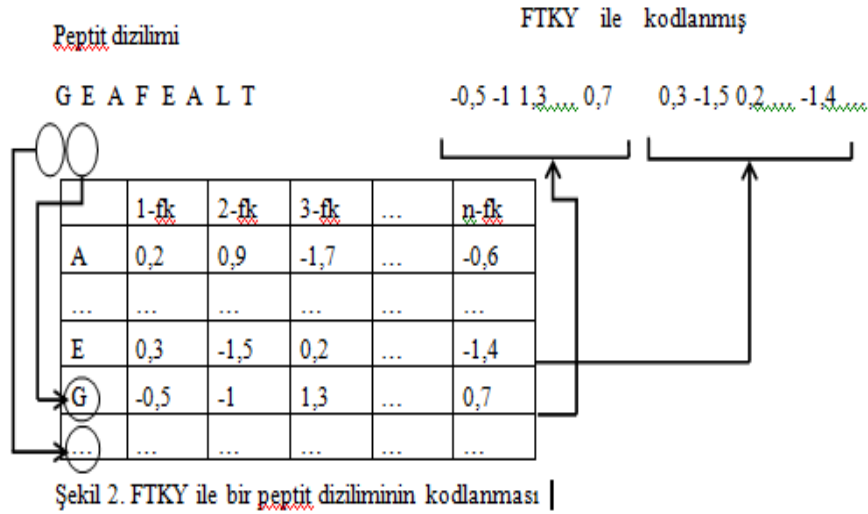
2.3. Destek Vektör Makineleri

DVM, 1979 yılında Vapnik tarafından geliştirilmiştir. DVM, eğitim örneklerinin, bir üst düzlem (hyperplane) ile doğrusal olarak ayrılabilirliği üzerine kurulu bir makine öğrenmesi yöntemidir. Eğitim için kullanılacak N örüntüden oluşan verinin, $D = \{x_i, y_i\}_{i=1,2,\dots,N}$, olduğunu varsayalım. Burada $x_i \in \mathcal{R}^d$ eğitim örnekleri ve $y_i \in \{-1, +1\}$ etiket değerleridir. Doğrusal olarak ayrılabilir durumda, iki sınıfa ayrılabilen örüntüler direkt olarak buldukları orijinal uzayda bir üst düzlem ile ayrılabilirler. DDVM'nin amacı ayırıcı üst düzlemin iki eğitim sınıfına eşit uzaklıkta olmasını sağlayarak eğitim örneklerini ayırmaktır. Eğer eğitim örüntüleri giriş uzayında doğrusal olarak bir üst düzlem ile ayrılmıyorlarsa, DVM bu eğitim örneklerine ait öznitelikler vektörlerini, yüksekboyutlu bir öznitelikler uzayına taşıyarak, bir üst düzlem ile doğrusal olarak iki sınıfa ayrılabilirliğini sağlamaktadır [13].

2.4. Birimlik Öznitelik Kodlama Yöntemi

HIV-1 proteaz enzimi bölünme kısımlarının tespitinde birimlik vektörlerle (birbirine dik birim vektör – orthonormal vector) kodlama yöntemi (BKY), en sık uygulanan yöntemlerdendir [14]. BKY, P peptitini oluşturan her bir P_i amino asit sembolü, birbirine dik, $d_i = (\delta_{i1}, \delta_{i2}, \dots, \delta_{i20})$ vektörlerle ifade edilir. Burada δ_{ij} Kronecker delta sembolüdür.

BKY'de her bir amino asit 20 bit uzunluğunda vektör ile temsil edilir. Bu temilde, her bir kalıntının sırasına karşılık gelen bit 1 ile, geri kalan değerler ise 0 ile temsil edilir [15]. Böylece her bir peptit dizilimi 1x160 bit büyüklüğünde birimlik vektörlerle temsil edilir.



FTKY, amino asitlerin fizikokimyasal özelliklerinin modellemeye yansıtılması üzerine kuruludur. Böylece HIV-1 proteaz kesim konumlarının tahmininde, peptitleri meydana getiren amino asitlerin birbirleri ile olan fizikokimyasal etkileşimleri sınıflandırıcıya daha iyi tanıtılmaktadır.

3. Bulgular

DeneySEL çalışmalarda, FTKY, 10-fk, 20-fk, 30-fk, 40-fk ve 50-fk'ye göre sınıf doğruluğu, duyarlık ve Alıcı İşletim Karakteristiği Eğrisi Altında Kalan Alan (AİKAA - AUROC) değerleri bakımından test edilmiştir. Sınıf doğruluğu değeri, modellemesi gerçekleştirilen HIV-1 proteaz enziminin kesme konumlarının tahmininde, doğru tahmin edilen peptit sayısının (kesme konumuna sahip olan ve olmayan), tüm peptit sayısına oranıdır [18]. Alıcı İşletim Karakteristiği (AİK - ROC) eğrisi, testin değişik kesim noktalarında doğru pozitif (y-ekseni) değerlerinin, yanlış pozitif değerlerine (x-ekseni) karşı noktalanması ile elde edilir. Her kesim noktasındaki doğru pozitif ve yanlış pozitive karşılık gelen noktalar birleştirilerek AİK eğrisi çizilir. AİKAA ise AİK eğrisi altında kalan alanın değeridir [19]. Duyarlık veya doğru pozitif değeri ise doğru tahmin edilen kesme konumuna sahip peptit sayısının, tüm kesme konumuna sahip peptit sayısına oranını ifade eder [18]. Geliştirilen yöntemin testleri, matematik ve grafik fonksiyonları üzerine kurulu, etkileşimli bir programlama ortamı olan MatLab programında gerçekleştirilmiştir. DDVM sınıflandırıcı yöntemi OSU Toolbox [20] ile uygulanmıştır.

Testler, 10-kat çapraz doğrulama tekniğine (ÇDT) göre gerçekleştirilmiştir. 10-kat ÇDT'de veri seti, 10 kümeye ayrılır. Kesilmiş peptitler ve kesilmemiş peptitler her bir kümeye rastgele ve eşit olacak şekilde dağıtılır. Bir çapraz doğrulamada, 10 kümeden 9'u eğitim verisi, 1'si test verisi olarak modelleme gerçekleştirilir. Bir testte toplam 10 çapraz doğrulama gerçekleştirilir [21]. Böylece her bir küme hem eğitim hem de test verisi olarak test sürecine dahil olur. Elde edilen sonuçlar 10 test üzerinden gerçekleştirilmiştir.

PR-1625 verileri üzerinde yapılan testlerde, Tablo 1'de görüldüğü gibi FTKY ile 10-fk'ye göre kodlanan girişler en yüksek sınıf doğruluğu değeri verirken 50-fk'da en düşük başarıyı vermiştir. Fizikokimyasal özellik sayısı arttıkça performans düşmektedir. FTKY hem doğrudan hem öznetelik çıkarım yöntemleri uygulandığında PR-3261 veri setinde PR-1625 veri setine göre daha düşük başarıyı sergilemiştir.

Tablo 1. FTKY'nin PR-1625 ve PR-3261 veri setleri üzerinde sınıf doğruluğu başarımı

	PR-1625 (%)	PR-3261 (%)
10-fk	95,21	93,11
20-fk	94,81	93,87
30-fk	94,78	94,37
40-fk	94,46	94,01
50-fk	94,19	94,03

Tablo 2'de ise FTKY'nin AİKAA değerleri görülmektedir. Bu sonuçlara göre PR-1625 üzerinde en

yüksek başarıyı 0,99 değeri ile 20-fk'ya göre yapılan kodlama vermiştir. Yine PR-1625 üzerinde yapılan test performansları PR-3261'e göre daha yüksektir. En düşük değer ise PR-3261 üzerinde, 10-fk'ya göre yapılan kodlamadan elde edilmiştir.

Tablo 2. FTKY'nin PR-1625 ve PR-3261 veri setleri üzerindeki karşılaştırmalı AİKAA sonuçları

	PR-1625	PR-3261
10-fk	0,98	0,95
20-fk	0,99	0,96
30-fk	0,98	0,96
40-fk	0,98	0,96
50-fk	0,98	0,96

Tablo 3'de FTKY'nin PR-1625 ve PR-3261 veri setleri üzerinde duyarlık değerleri görülmektedir. Yapılan testlerde en yüksek sonuçlar yine PR-1625 üzerinde elde edilmiştir. En yüksek duyarlık değeri PR-1625'de % 90,59 ile 30-fk ile yapılan kodlamada, PR-3261'de ise % 76,60 değeri ile 50-fk'ya göre yapılan kodlamada elde edilmiştir.

Tablo 3. PR-1625 ve PR-3261 veri setleri üzerinde FTKY'nin duyarlık başarımları

	PR-1625 (%)	PR-3261 (%)
10-fk	89,3	70,32
20-fk	90,11	74,51
30-fk	90,59	76,18
40-fk	89,95	75,48
50-fk	89,65	76,6

4. Tartışma ve Sonuç

Bu çalışmada HIV-1 enzimi bölünme kısımlarının tespiti için geliştirilen FTKY, DDVM yöntemi kullanılarak güncel PR-1625 ve PR-3261 veri setleri üzerinde test edilmiştir. Buna göre FTKY, PR-1625 üzerinde en iyi 10-fk özelliğe göre, PR-3261 üzerinde en iyi 30-fk'ya göre yapılan kodlamalarda diğerlerine nazaran yüksek sınıf doğruluğu başarımları sergilemiştir. FTKY'nin amino asitlerin fizikokimyasal özellikleri göz önüne alarak geliştirilmiş olması, örüntülerin sınıflandırıcıya daha iyi tanıtılmasının yolunu açmıştır. Bu durum tek bir sınıflandırıcı üzerinde bile yüksek sınıf doğruluğu ve AİKAA değerleri elde edilmesinden anlaşılmaktadır. FTKY, PR-1625 ve PR-3261 veri setleri üzerinde doğrusal lojistik (linear logistic), doğrusal perseptron (linear perceptron), doğrusal ayırıcı (linear discriminant)

sınıflandırıcıları ile quadratik (quadratic), parzen, naive bayes doğrusal olmayan sınıflandırıcıları ile de test edilmiştir. Fakat bu algoritmaların hiçbirisi DDVM kadar başarılı sonuçlar verememişlerdir. PR-1625 üzerinde gerçekleştirilen testlerde PR-3261'e göre daha yüksek performanslar elde edilmiştir. Bu durum PR-3261 veri setinin varyansının daha yüksek olduğunu göstermektedir. FTKY'nin üzerine gelecekte yapılacak çalışmalarda en iyi özelliklerin seçiminde fizikokimyasal özelliklerin birbirinden bağımsız olarak düşünülmemesi başarımları artırılabilir. Bu çerçevede FTKY'de birinci safhada en iyi özellikler seçilirken amino asitler arasındaki bağımlılığın kodlamaya dahil edilebilmesi için,

$$\binom{544}{n}, \quad (1)$$

kadar seçenek denenebilir. (1)'deki notasyonda 544 fizikokimyasal özelliğin n 'li kombinasyonları hesaplanmaktadır. Burada n , en iyi fizikokimyasal özellik sayısıdır. Bu işlem sonunda en az sınıflandırma hatası yapan n adet fizikokimyasal özellik belirlenebilir. Böylece daha yüksek başarımlar elde edilebilir. Ayrıca FTKY'nin, diğer makine öğrenmesi yöntemleri ile başka sınıflandırma problemlerinin çözümünde de uygulanması gerçekleştirilebilir.

Kaynaklar

1. Beck Z.Q., Hervio L, Dawson P.E., Elder J.H., Madison E.L., Identification of efficiently cleaved substrates for HIV-1 protease using a phage display library and use in inhibitor development. *Virology* 274 (2):391-401, 2000.
2. Graves, B.J., Hatada, M.H., Miller, J. K., Graves, M.C., Roy, S., Cook, C.M., Krohn, A., Martin, J.A., Roberts, N.A., In *Structure and Function of the Aspartic Protease: Genetics, Structure and Mechanisms*. Dunn, B., Ed. Plenum: New York; p. 455, 1992.
3. Kuo-Chen Chou, Prediction of Human Immunodeficiency Virus Protease Cleavage Sites in Proteins, *Analytical Biochemistry*. 233, 1–14, 1996.
4. [4] Cai Y.D., Chou K.C., Artificial neural network model for predicting HIV protease cleavage sites in protein. *Adv Eng Software*, 29:119–128, 1998.
5. Rönngvaldsson, T., You, L., Why neural networks should not be used for HIV-1 protease cleavage site prediction. *Bioinformatics*, 1702–1709, 2003.
6. Nanni L, Lumini A., MppS: an ensemble of Support Vector Machine based on multiple physicochemical properties of amino-acids. *Neuro Computing*, 69:1688-1690, 2006.

7. Carlotta Orsenigo, Carlo Vercellis, Predicting HIV Protease-Cleavable Peptides by Discrete Support Vector Machines. *Machine Learning and Data Mining in Bioinformatics*, 197-206, 2007.
8. Kawashima, S., Kanehisa, M., AAindex: amino acid index database, *Nucleic Acids Res.* 20 (1): 374, 2000. (www.genome.jp/aaindex/)
9. Schechter, I., Berger, A. 1967. On the size of the active site in proteases. *Biochemical and Biophysical Research Communications* 27, 157–162, 1967.
10. Kontijevskis, A., Wikberg, J.E., Computational proteomics analysis of HIV-1 protease interactome. *Proteins-Structure Function and Bioinformatics* 68(1): pp. 305-312, 2007.
11. Schilling, O., Overall, C.M., Proteome-derived, database-searchable peptide libraries for identifying protease cleavage sites. *Nat Biotechnol* 26(6): pp. 685-694, 2008.
12. Rognvaldsson, T., Etchells, T.A., How to find simple and accurate rules for viral protease cleavage specificities. *BMC Bioinformatics* 10: 149, 2009.
13. Wang, L., *Support Vector Machines: Theory and Applications*, Springer, 2005.
14. Gök M., Özcerit A.T., Linear Support Vector Machines for HIV-1 Protease Site Detection, *ISSD'09*, Sarajevo, Bosnia Herzegovia, pp. 381-384, 2009.
15. [15] Narayanan, A., Wu, X., Yang, Z.R., Mining viral protease data to extract cleavage knowledge. *Bioinformatics*, 18: Suppl 1, pp. 5-13, 2002.
16. Akdemir, B., Tahmin uygulamalarında performans geliştirmek için kullanılan normalizasyon metotlarına yeni bir yaklaşım. Doktora Tezi, Selçuk Üniversitesi, 2009.
17. Jain A., Nandakumar K., Score normalization in multimodal biometric systems. *Pattern Recognition*, 38(12): pp. 2270-2285, 2005.
18. Fawcett, T., ROC graphs: Notes and practical considerations for researchers. Technical Report, HP Laboratories. California, USA, 2004.
19. Elmali, F., Altın Standartlı ve Altın Standartsız Durumlarda, Yarı Parametrik ve Parametrik Olmayan ROC eğrisi Yöntemlerinin Karşılaştırılması. Osmangazi Üniversitesi, Sağlık Bilimleri Enstitüsü, Doktora Tezi, 2009.
20. Junshui, M.A., Y.I., Zhao., *OSU SVM Toolbox for MATLAB*, 2002. (<http://sourceforge.net/projects/svm/>)
21. Duda, R.O., Hart, P.E., Stork, D.G., *Pattern Classification*, 2nd edition. John Wiley & Sons Inc, 2001