

ヒト免疫不全ウイルスに対する抗ウイルス薬の  
選択に関する in silico 研究

2020

太田 亮作



# 目次

|                                      |    |
|--------------------------------------|----|
| 総論の部                                 | 1  |
| 緒言                                   | 1  |
| <br>                                 |    |
| 第一章                                  |    |
| モデルベースメタ解析法による初回療法での最も有効なレジメンの選<br>択 | 2  |
| I-1 データの収集                           | 3  |
| I-2 モデルの構築および最適化                     | 7  |
| I-3 共変量モデル解析                         | 11 |
| I-4 考察                               | 15 |
| <br>                                 |    |
| 第二章                                  |    |
| HIVの薬剤耐性に関する予測モデル構築法の提案と検証           | 17 |
| II - 1 データ収集とデータの前処理                 | 19 |
| II - 2 分子場ポテンシャルの計算                  | 22 |
| II - 3 予測モデルの構築                      | 22 |
| II - 3 - a 特徴抽出 . . . . .            | 23 |
| II - 3 - b ハイパーパラメータの最適化 . . . . .   | 25 |
| II - 3 - c モデルの最終評価 . . . . .        | 27 |
| II - 4 薬剤耐性ウイルスの構造要因解析               | 30 |
| II - 5 考察                            | 32 |
| <br>                                 |    |
| 結論                                   | 34 |

|          |    |
|----------|----|
| 謝辞       | 36 |
| 実験の部     | 37 |
| 第一章 実験の部 | 37 |
| 第二章 実験の部 | 39 |
| 引用文献     | 40 |

# 総論の部

## 緒言

ヒト免疫不全ウイルス (human immunodeficiency virus : HIV) の感染は、臨床上大きな問題である。HIV感染者は、現在約3800万人と推定されており [1]、エイズ関連疾患による死者は毎年約70万人にのぼる [1]。世界保健機関 (World Health Organization:WHO) を含む国際機関が協力し構築された国連合同エイズ計画 (Joint United Nations Programme on HIV/AIDS:UNAIDS) では、HIVに関して、2020年までに新規感染者およびエイズ関連疾患による死者をそれぞれ50万人以下に抑えること、さらに、90%の感染者が自身の感染を知り、認識した感染者の90%が治療を行い、治療を選択した感染者の90%で治療を成功させることを目標 (90-90-90 goals) として定めた [2]。現状ではこれらの目標の達成が難しいとされるものの、医療環境の整備が非先進国でも進んだ結果、薬剤の供給が問題であった過去とは異なり、現在では自由に治療薬を選ぶことが国によらず可能となりつつある。

現状の抗ウイルス治療は、少なくとも感染者の10%で十分な治療効果が得られていない [2]。不十分な治療は、エイズの発症リスクを高めるとともに、薬剤耐性ウイルスの発生確率を高めるといった報告もある [3]。社会的には、不十分な治療が薬剤耐性ウイルスを増やし、さらに治療効果が低下するといった悪循環をもたらす可能性がある。治療効果の改善の抜本的解決法には、新規薬剤の開発が望まれるが、薬剤耐性を獲得しやすいというHIVの性質上、開発頻度を高くする必要があり、開発のハードルも着実に高まっていくことから、社会への多大な経済的負担が課題となる。

こういった背景から、著者は、既存の薬剤の効率的な利用が重要と考え、その治療効果を社会レベルで最大化するためには、感染者の大多数に適用されることとなる第一選択レジメンを最適化しておくことおよび少数の薬剤耐性ウイルス感染者への治療戦略の確立というマクロ・ミクロの両面からの研究が必要と考えた。前者においては、数理モデルを用いることで、従来不可能とされてきた質的に異なる臨床試験の結果を統合し、時間に対する定量的な解釈を可能とするモデルベースメタ解析法を導入し、レジメンを大規模かつ多角的に評価した。後者においては、ウイルスの薬剤耐性に基づく個別的治疗の実現に向けて、感染したウイルスの耐性を高い精度で予測することを期待し、タンパク質の立体構造に焦点をあてて、その構造に基づく生理学的・物理化学的情報を組み込んだ機械学習モデルの構築を試みた。

以下、二章にわたり本研究で得られた知見について論述する。

## 第一章

# モデルベースメタ解析法による初回療法での最も有効なレジメンの選択

抗HIV治療には、Nucleoside Analogue Reverse Transcriptase Inhibitor (NRTI) と Non-Nucleoside Reverse Transcriptase Inhibitor (NNRTI) から二種類のバックボーンと、Integrase inhibitor (INSTI)、Protease inhibitor (PI) のいずれか一種類のキードラッグからなる多剤併用療法が採用されている。長期的なウイルス抑制には、初回療法の選択が重要であることが知られており [4]、NNRTIであるエファビレンツは第一選択のキードラッグとして長年使用されてきた [5, 6]。しかし、INSTIの一つであるドルテグラビルが2013年に承認され、2年も経たないうちに、Department of Health and Human Services (DHHS) と European AIDS Clinical Society (EACS) がドルテグラビルを第一選択のキードラッグとして推奨するようになった [7]。その結果、エファビレンツは、先進国ではほとんど使用されなくなったが、医療経済学的な理由により、発展途上国を含む世界全体で見れば未だに多くのエファビレンツが使用されているという現状がある [7]。実際、世界保健機関 (WHO) がキードラッグとしてエファビレンツよりもドルテグラビルを優先的に使用するよう提唱するようになったのは、当該研究の実施後のことである [8]。

ドルテグラビルのキードラッグとしての有効性は、システマティックレビューやその解析によって実証されてきた。Rutherfordら [9] は、抗HIV治療未経験の成人患者を対象に、ドルテグラビルベースのレジメンとエファビレンツベースのレジメンの有効性を比較するメタ解析を行った。その結果、ドルテグラビルベースのレジメンが、持続的なウイルス抑制、薬剤耐性の出現率の低さ、免疫の回復の早さといった点で優れていることが示された。しかし、Rutherford [9] らのメタアナリシスは2つの臨床試験のみを対象とした小規模なものに過ぎず、結論を導くのに十分な評価がなされたとは言い難い。

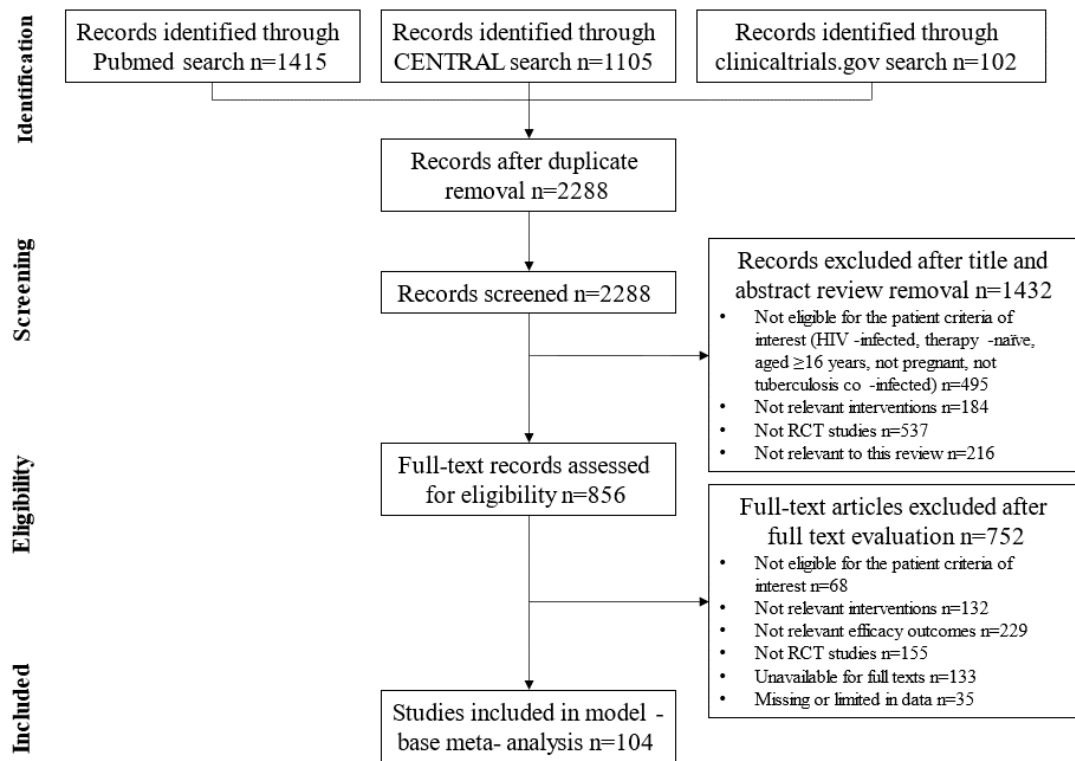
こうした背景の中で、質的に異なるデータを統合的に扱うことができる Model-based meta-analysis (MBMA) [10, 11] を解析に用いることとした。MBMAは、治療効果を時間の関数としてモデル化することで、投与量、期間、サンプリング時間の異なる臨床試験の情報を統合し、定量的な解釈を可能とする。本章では、MBMAを用いて、既存の研究手法では統合できなかった複数の臨床試験のデータから、ドルテグラビル、エファビレンツの各レジメンの有効性の比較を行った。

## I-1 データの収集

臨床試験データ収集は、Pubmed、Cochrane Central Register of Controlled Trials (CENTRAL)、clinicaltrials.gov の複数データベースから取得した。試験の論文を取得する検索ワードは、ドルテグラビルには Rutherford ら [9] の論文、エファビレンツには Kryst らの論文 [12] で採用されているワードを利用した。適格基準の設定は、コクランハンドブックに基づいて行った。コクランハンドブック [13] は、収集するデータを均質化するためのコンセプトを定めたものであり、本研究では、初回療法のレジメンの比較を前提に、以下の適格基準を具体的に定めた。

- 臨床試験のデザインは、ランダム化比較試験であること
- 1 アーム 10 人以下の臨床試験を除くこと
- 患者は、過去に治療経験のない 16 歳以上の HIV 感染者であること。ただし、結核を併発している患者や妊娠中の患者を除くこと
- 対象となる初回療法のレジメンは、dolutegravir (DTG) + abacavir (ABC)+lamivudine (3TC), efavirenz (EFV)+ABC+3TC, EFV + tenofovir disoproxil fumarate (TDF) + emtricitabine (FTC)、EFV + TDF + 3TC のいずれかであること
- 治療効果は、血中のウイルス濃度が 50 copies/mL 以下の患者の割合で評価されていること

これらの適格基準に合う論文を、PRISMA ガイドライン [14] に従い、Figure.1 に示す手順で抽出した。その結果、Pubmed、CENTRAL、clinicaltrials.gov から集められた計 2622 報の論文から、最終的に 104 報の論文が選抜された。同じ臨床試験を指しているものの複数に分散して投稿された文献の重複を削除し整理した結果、Table 1 に示す計 30 の臨床試験データが得られた [15]{45}。



**Figure 1.** PRISMA statement 2009 flow diagram. The diagram depicts the selection process of studies undertaken in the present meta-analysis. Reasons for exclusion are provided along with their relevant counts.



**Table 1.** List of data used for the present analysis

| Source                    | Study                | No. of patients | Regimen     | Time points  | Proportion of participants with viral load of < 50 copies/ml                                   | HIV-1 RNA log <sub>10</sub> copies/ml | CD4 cells=mm <sup>3</sup> |
|---------------------------|----------------------|-----------------|-------------|--|--|---------------------------------------|---------------------------|
| Walmsley S. et al. [15]   | SINGLE               | 419             | EFV+ABC+3TC | 4, 8, 12, 16, 24, 32, 40, 48, 60, 72, 84, 96                 | 0.14, 0.35, 0.55, 0.69, 0.84, 0.82, 0.80, 0.78, 0.76, 0.75, 0.72                               | 4.68                                  | 338                       |
| Sax PE. et al. [16]       | A5202                | 388             | EFV+ABC+3TC | 16, 24, 36, 48, 60, 72, 84, 96                               | 0.43, 0.60, 0.72, 0.74, 0.78, 0.73, 0.76, 0.85   | 5                                     | 138                       |
|                           |                      | 399             | EFV+TDF+FTC | 16, 24, 36, 48, 60, 72, 84, 96                               | 0.46, 0.69, 0.79, 0.80, 0.80, 0.84, 0.79, 0.84   | 5                                     | 146                       |
| Honda M. et al. [17]      | NCT00280969          | 36              | EFV+ABC+3TC | 4, 8, 12, 16, 20, 24, 28, 32, 36, 40, 44, 48, 60, 72, 84, 96 | 0.17, 0.32, 0.53, 0.68, 0.78, 0.78, 0.83, 0.78, 0.75, 0.84, 0.78, 0.78, 0.75, 0.76, 0.69, 0.56 | 4.6                                   | 220                       |
| Echeverria P et al. [18]  | NCT0031812           | 63              | EFV+ABC+3TC | 48   | 0.57   | 5.40 <sup>a</sup>                     | 193 <sup>a</sup>          |
| Post FA. et al. [19]      | ASSERT (NCT00549198) | 192             | EFV+ABC+3TC | 4, 12, 24, 36, 48  | 0.12, 0.50, 0.65, 0.69, 0.59   | 5.01                                  | 240                       |
|                           |                      | 193             | EFV+TDF+FTC | 4, 12, 24, 36, 48  | 0.087, 0.49, 0.75, 0.80, 0.70  | 5.12                                  | 230                       |
| Kumar P. et al. [20]      | SUPPORT              | 50              | EFV+ABC+3TC | 2, 4, 8, 12, 24, 36, 48, 60, 72, 84, 96                      | 0.066, 0.27, 0.50, 0.71, 0.87, 0.81, 0.81, 0.77, 0.67, 0.71, 0.67                              | 4.82                                  | 272.5                     |
| Bartlett JA. et al. [21]  | CLASS                | 97              | EFV+ABC+3TC | 8, 16, 24, 32, 40, 48, 56, 64, 72, 80, 88, 96                | 0.33, 0.63, 0.78, 0.7031, 0.76, 0.75, 0.73, 0.73, 0.73, 0.72, 0.70, 0.68                       | 4.9                                   | 307                       |
| Moyle G.J. et al. [22]    | CNA30021             | 384             | EFV+ABC+3TC | 48   | 0.66   | 4.91                                  | 264                       |
|                           |                      | 386             | EFV+ABC+3TC | 48   | 0.68   | 4.87                                  | 259                       |
| DeJesus E. et al. [23]    | CNA30024             | 324             | EFV+ABC+3TC | 2, 4, 8, 12, 16, 24, 36, 48                                  | 0.054, 0.11, 0.27, 0.46, 0.60, 0.71, 0.70, 0.71  | 4.81                                  | 267                       |
| Podzamczar D. et al. [24] | ABCDE                | 115             | EFV+ABC+3TC | 96   | 0.61   | 4.96                                  | 175                       |
| Kravchenko A. et al [25]  | NA                   | 60              | EFV+TDF+FTC | 24, 48   | 0.67, 0.84   | 4.8                                   | 310                       |
| Arribas JR. et al. [26]   | GS-01-934            | 244             | EFV+TDF+FTC | 16, 24, 32, 48, 72, 96                                       | 0.65, 0.78, 0.81, 0.79, 0.71, 0.700, 0.65, 0.64  | 5                                     | 237                       |
| Rockstroh JK. et al. [27] | STARTMRK             | 282             | EFV+TDF+FTC | 2, 4, 8, 12, 16, 24, 32, 40, 48, 60, 72, 84, 96              | 0.024, 0.12, 0.39, 0.61, 0.79, 0.85, 0.85, 0.83, 0.83, 0.80, 0.83, 0.80, 0.79                  | 5                                     | 217.4                     |
| Molina JM. et al. [28]    | ECHO                 | 344             | EFV+TDF+FTC | 2, 4, 8, 12, 16, 24, 32, 40, 48                              | 0.045, 0.14, 0.35, 0.58, 0.77, 0.84, 0.84, 0.83, 0.82  | 5                                     | 257                       |
| Landman R. et al. [29]    | DAYANA               | 30              | EFV+TDF+FTC | 16, 24, 36, 48, 72, 96                                       | 0.53, 0.53, 0.73, 0.73, 0.747, 0.73  | 5.6                                   | 201                       |

| Source                    | Study                        | No. of patients | Regimen     | Time points                                     | Proportion of participants with viral load of < 50 copies/mL                 | HIV-1 RNA log <sub>10</sub> copies/ml | CD4 cells=mm <sup>3</sup> |
|---------------------------|------------------------------|-----------------|-------------|---|--|---------------------------------------|---------------------------|
| Vernazza P. et al. [30]   | A5271015                     | 63              | EFV+TDF+FTC | 2, 4, 8, 12, 16, 24, 32, 40, 48                 | 0.13, 0.30, 0.48, 0.75, 0.78, 0.88, 0.88, 0.87, 0.86                         | 4.7                                   | 312                       |
| Cohen C. et al. [31]      | NCT00869557                  | 23              | EFV+TDF+FTC | 2, 4, 8, 12, 16, 24, 32, 40, 48                 | 0.086, 0.25, 0.53, 0.74, 0.83, 0.82, 0.87, 0.88, 0.82                        | 4.58 <sup>a</sup>                     | 436                       |
| Amin J. et al. [32]       | ENCORE1                      | 309             | EFV+TDF+FTC | 48  | 0.78   | 4.73 <sup>a</sup>                     | 272 <sup>a</sup>          |
| Amin J. et al. [33]       | ENCORE1                      | 309             | EFV+TDF+FTC | 96  | 0.8  | 4.73 <sup>a</sup>                     | 272 <sup>a</sup>          |
| Sax PE et al. [34]        | GS-US-236-0102 (NCT01095796) | 352             | EFV+TDF+FTC | 2, 4, 8, 12, 16, 24, 32, 40, 48                 | 0.074, 0.20, 0.50, 0.68, 0.81, 0.86, 0.86, 0.85, 0.86                        | 4.78                                  | 383                       |
| Van Lunzen J. et al. [35] | STaR (NCT01309243)           | 392             | EFV+TDF+FTC | 4, 8, 12, 16, 24, 32, 40, 48, 60, 72, 84, 96    | 0.48, 0.44, 0.66, 0.79, 0.87, 0.86, 0.87, 0.86, 0.85, 0.84, 0.81, 0.78       | 4.80 <sup>a</sup>                     | 385.2 <sup>a</sup>        |
| Thompson M. et al. [36]   | 652-2-202 (NCT01338883)      | 28              | EFV+TDF+FTC | 24, 48  | 0.71, 0.50   | 4.56                                  | 310                       |
| Miro JM. et al. [37]      | Advanz-3                     | 28              | EFV+TDF+FTC | 12, 24, 36, 48                                  | 0.43, 0.71, 0.75, 0.64   | 5.12                                  | 41                        |
| Puls, RL. et al. [38]     | Altair                       | 114             | EFV+TDF+FTC | 48  | 0.85   | 4.67                                  | 227                       |
| Markowitz M. et al. [39]  | Protocol 004 (NCT00100048)   | 38              | EFV+TDF+3TC | 2, 4, 8, 12, 16, 24, 32, 40, 48, 60, 72, 84, 96 | 0.11, 0.23, 0.38, 0.68, 0.71, 0.93, 0.90, 0.87, 0.87, 0.84, 0.87, 0.82, 0.84 | 4.88 <sup>b</sup>                     | 276 <sup>a</sup>          |
| Gugliotti R et al. [40]   | NCT00350272                  | 37              | EFV+TDF+3TC | 2, 4, 6, 8, 10, 12                              | 0.080, 0.11, 0.35, 0.54, 0.57, 0.70  | NA                                    | NA                        |
| Gallant JE. et al [41]    | Study 903                    | 299             | EFV+TDF+3TC | 48, 96, 144                                     | 0.76, 0.73, 0.68   | 4.91 <sup>a</sup>                     | 276 <sup>a</sup>          |
| Ra F. et al. [42,43]      | SPRING-2                     | 411             | DTG+ABC+3TC | 4, 8, 12, 16, 24, 32, 40, 48, 60, 72, 84, 96    | 0.71, 0.85, 0.87, 0.89, 0.93, 0.91, 0.88, 0.88, 0.87, 0.86, 0.84, 0.81       | 4.52                                  | 359                       |
| Gallant J. et al. [44]    | GS-US-380-1489               | 315             | DTG+ABC+3TC | 48  | 0.93   | 4.42                                  | 443                       |
| Orrell C. et al. [45]     | ARIA                         | 248             | DTG+ABC+3TC | 4,12,24,36,48                                   | 0.93, 0.64, 0.81, 0.85, 0.85, 0.82   | 4.41                                  | 340                       |

a .mean, b :geometric mean, NA:Not Applicable  
EFV+ABC+3TC means each dose is 600mg, 600mg, and 300mg. EFV+TDF+FTC means each dose is 600mg, 300mg, and 200mg.  
EFV+TDF+3TC means each dose is 600mg, 300mg, and 300mg. DTG+ABC+3TC means each dose is 50mg, 600mg, and 300mg

## I-2 モデルの構築および最適化

治療効果は、逐次一次反応に従って経時的に上昇し、減衰すると仮定して、以下の式で記述し、非線形混合効果モデルで解析を行った。

$$E = \bar{E} + \frac{S}{N} \frac{\bar{E}(1 - \bar{E})}{N} \quad (1)$$

$$\bar{E} = E_{max} (1 - e^{-k_o(t - t_{lag})}) e^{-k_d(t - t_{lag})} \quad (2)$$

$$k_o = \bar{k}_o e^{k_o} \quad (3)$$

$$k_d = \bar{k}_d e^{k_d} \quad (4)$$

$E$  : Drug effect

$E_{max}$  : Maximum drug effect

$k_o$  : Rate constant describing onset of drug effect

$k_d$  : Rate constant describing decay of drug effect

$N$  : Number of patients in the study

$t_{lag}$  : Time lag of onset of drug effect

$k_o, k_d$  : Inter study variability

$S$  : Random residual error

解析の結果推定された各レジメンのパラメータはTable 2の通りである。モデルによって推定された平均値のプロファイルを図示すると、収集されたタイムコースデータとよくフィットしていた (Figure 2)。さらに、推定パラメータの computer-calculated standard error に基づき、モンテカルロ法 (1000 回) を利用して、タイムコースの 90% 信頼区間を求めた。実測のデータは、ほとんど 90% 信頼区間に含まれ、平均値を中心として均等に分散しており、シミュレーションベースのモデルの予測性能の評価法の一つである視覚的事後予測性能評価 (visual predictive check) の結果は良好であった。

Table 2 の各レジメンの推定パラメータを比較することで、レジメンの効果を定量的に評価することができる。最大の治療効果を表す  $E_{max}$  は、各レジメンについてほとんど差がなかったが、治療効果発現速度を表す  $k_o$  は約 4 倍、ドルテグラビルベースのレジメンがエファビレンツベースのレジメンよりも大きかった。一方、治療効果減衰速度を表す  $k_d$  は、約 2 倍、エファビレンツベースのレジメンが大きかった。これらのパ

ラメータの比較は、ドルテグラビルベースのレジメンが治療効果の発現が早く、治療効果の減衰が遅いことを示している。

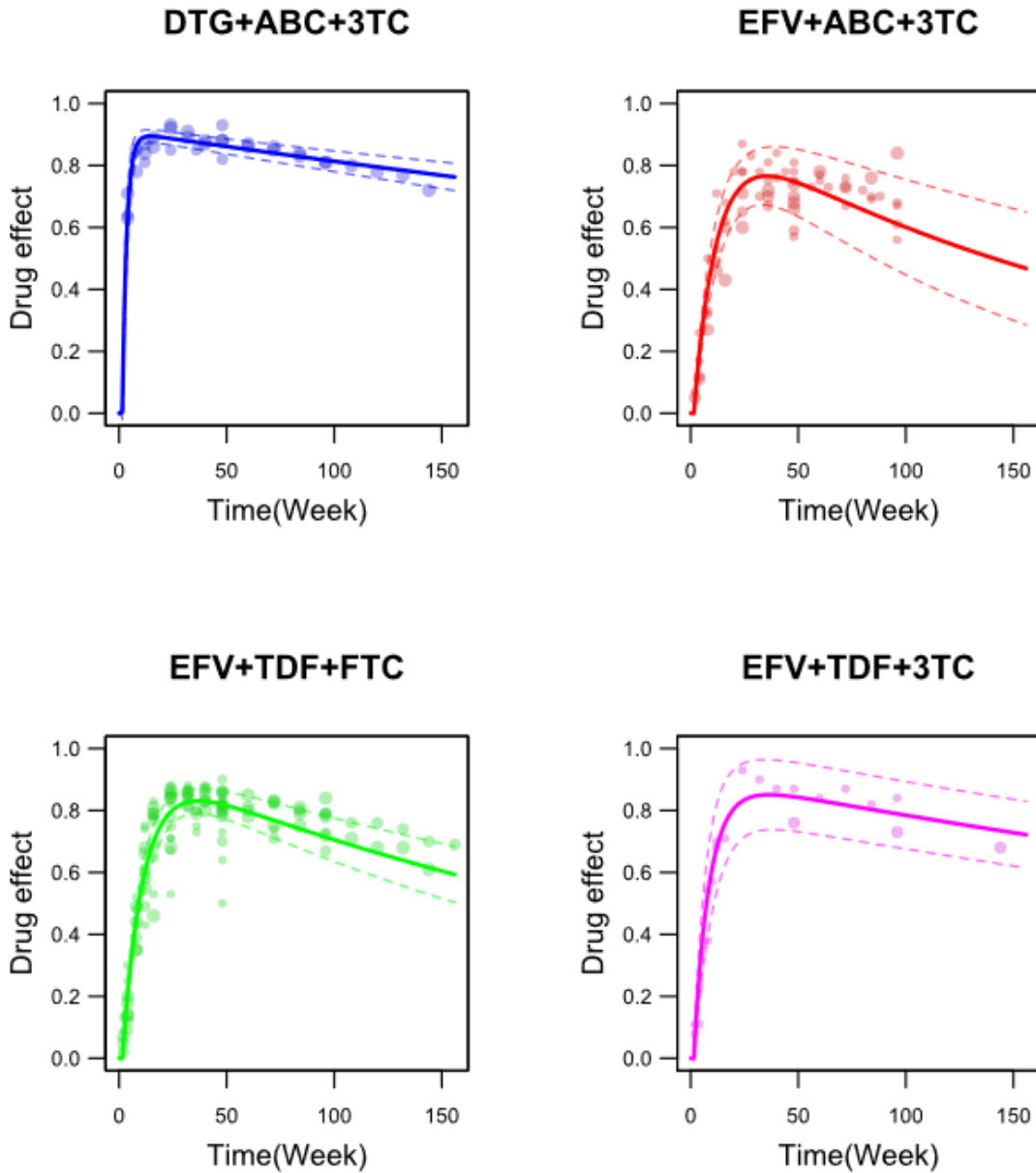
診断プロットを用いて、モデルの推定精度を評価した (Figure 3)。Goodness-of-fit plot の結果より、自由度調整済み決定係数が 0.89 と良好な適合度であることが分かった。また、時間に対する条件付重みつき残差 (conditional weighted residual, CWRES) は、ゼロを中心に、対称的に分布してかつレンジが  $\pm 3$  以内にほとんど収まっており、良好な結果であった。

**Table 2.** Parameter estimates for each regimen based on the base population model.

| Parameters          | Parameter estimates <sup>b</sup> |                     |                     |                   |
|---------------------|----------------------------------|---------------------|---------------------|-------------------|
|                     | DTG+ABC+3TC                      | EFV+ABC+3TC         | EFV+TDF+FTC         | EFV+TDF+3TC       |
| $E_{max}$           | 0.833 $\pm$ 0.014                | 0.905 $\pm$ 0.052   | 0.963 $\pm$ 0.015   | 0.910 $\pm$ 0.076 |
| $k_o$               | 0.474 $\pm$ 0.052                | 0.0779 $\pm$ 0.0087 | 0.0862 $\pm$ 0.0076 | 0.118 $\pm$ 0.021 |
| $k_o^a$             | 0.0720 $\pm$ 0.0274              |                     |                     |                   |
| $k_d$ ( $10^{-3}$ ) | 0.864 $\pm$ 0.356                | 3.83 $\pm$ 1.41     | 2.81 $\pm$ 0.47     | 1.72 $\pm$ 0.26   |
| $k_d^a$             | 0.149 $\pm$ 0.018                |                     |                     |                   |
| $t_{lag}^a$         | 1.55 $\pm$ 0.07                  |                     |                     |                   |
| $a$                 | 1.95 $\pm$ 0.36                  |                     |                     |                   |

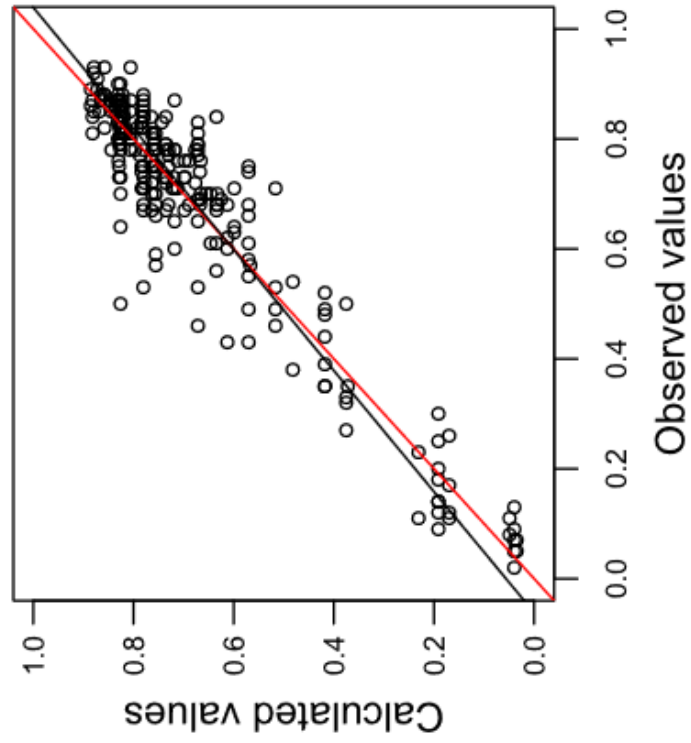
<sup>a</sup>  $t_{lag}$ , inter-individual variation of  $k_o$  and  $k_d$  ( $k_o$  and  $k_d$ ), and intra-individual variation ( ) were assumed to be constant irrespective to any of regimens, in order to in ation of number of parameters to be estimated by curve-fit.

<sup>b</sup> mean computer-calculated standard error

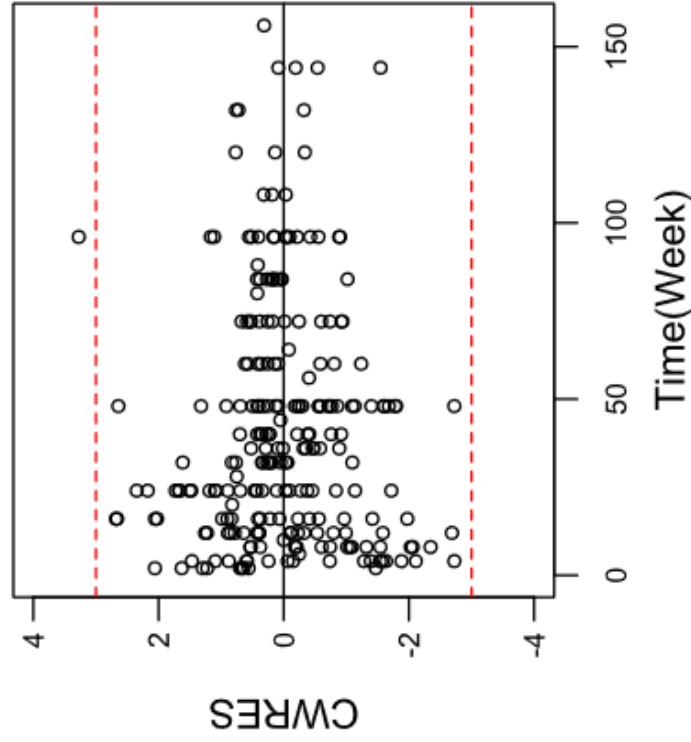


**Figure 2.** Time courses of observed and simulated drug effect for each regimen. The drug effect was defined as proportion of participants with viral load of < 50 copies/mL in plasma. Symbols represent observed data, of which the size is proportional to the number of patients in the studies. Each solid line represents the median, and the corresponding dashed lines represent the 5th and 95th percentiles estimated by exploiting 1000 times Monte Carlo simulations. Parameters used for the simulation are summarized in Table 2.

**Goodness of fit plot**



**CWRES against time plot**



**Figure 3.** Goodness of fit and CWRES plots associated with the base population model. Black and red lines represent regression and 1:1 correspondence lines, respectively.

### I-3 共変量モデル解析

治療効果の速度論パラメータと治療開始前の患者固有のパラメータとの関係を調べるために、治療開始直前の体内ウイルス濃度 (BSL)、治療開始直前の正常 CD4 のカウント数 (CD4)、 $\log_{10}k_o$ 、 $\log_{10}k_d$  の 4 変数の散布図行列を作成した (Figure 4)。その結果、 $\log_{10}k_o$  は、BSL と弱く正の相関を持っており、CD4 とは、弱く負の相関を持っていた。同様に、 $\log_{10}k_d$  も、BSL と弱く正の相関を持っており、CD4 とは、弱く負の相関を持っていた。このことは、重篤な患者ほど、治療効果の発現が遅く、治療効果の減衰が早い、つまり、予後が悪いということを示している。さらに、レジメンの違いの影響を Figure 4 の散布図行列から除いて解析するために、層別化回帰分析を行った結果 (Figure 5, Table 3)、散布図行列から得られた傾向がより顕著に表れた。

先にあげた患者固有のパラメータを共変量としてモデルへ組み込むことで、アーム間での患者情報の違いを補正することが可能となる。そこで、BSL もしくは CD4 を共変量として下のモデル式を仮定し、先の速度論モデルに組み込んだ。

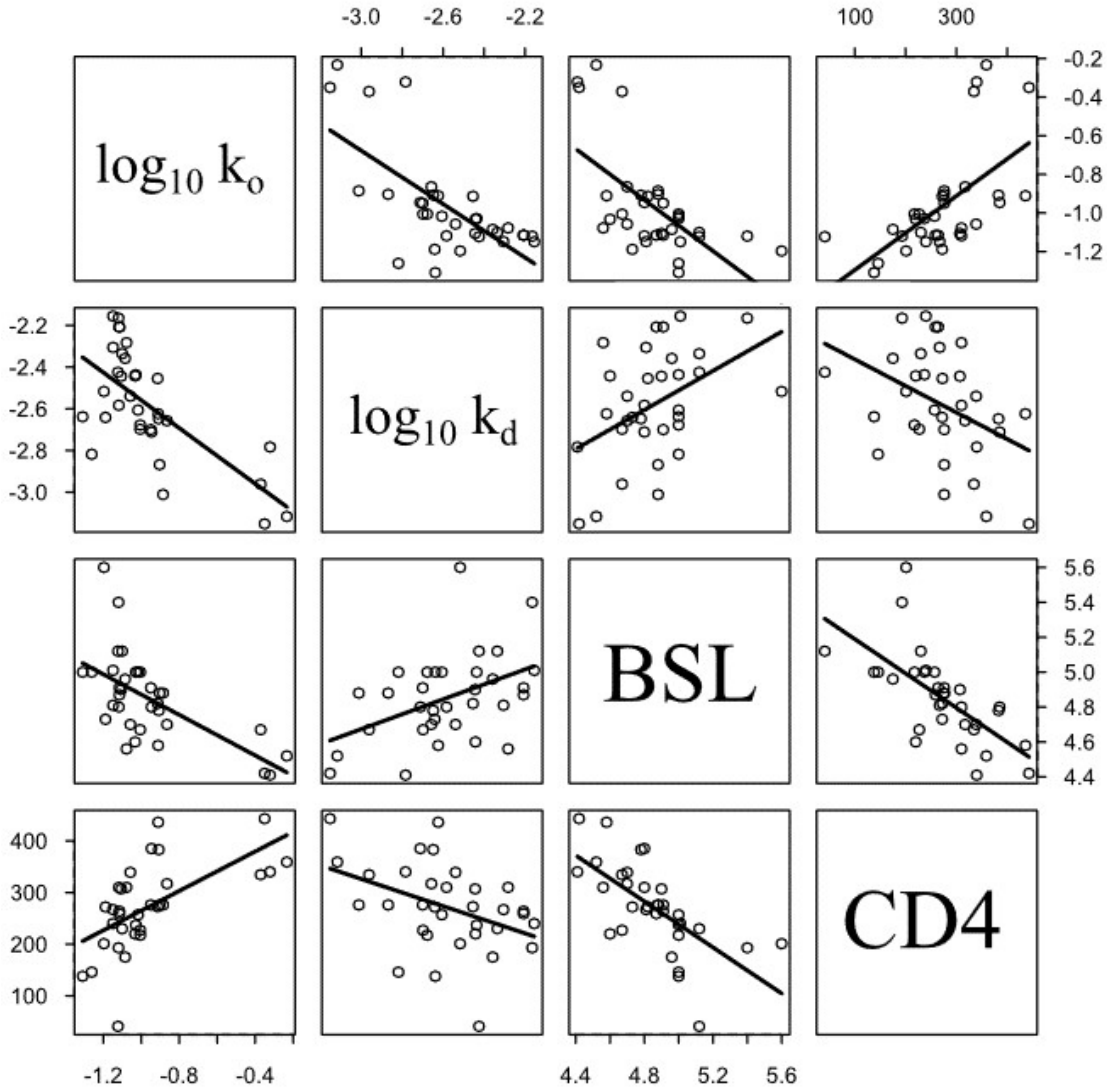
$$k_{o,i} = \theta_1 + \theta_2 \frac{COV_i - COV_{med}}{COV_{med}} \quad (5)$$

$k_{o,i}$  : Covariate-adjusted mean rate constant for study arm  $i$

$\theta_1, \theta_2$  : Estimated parameters

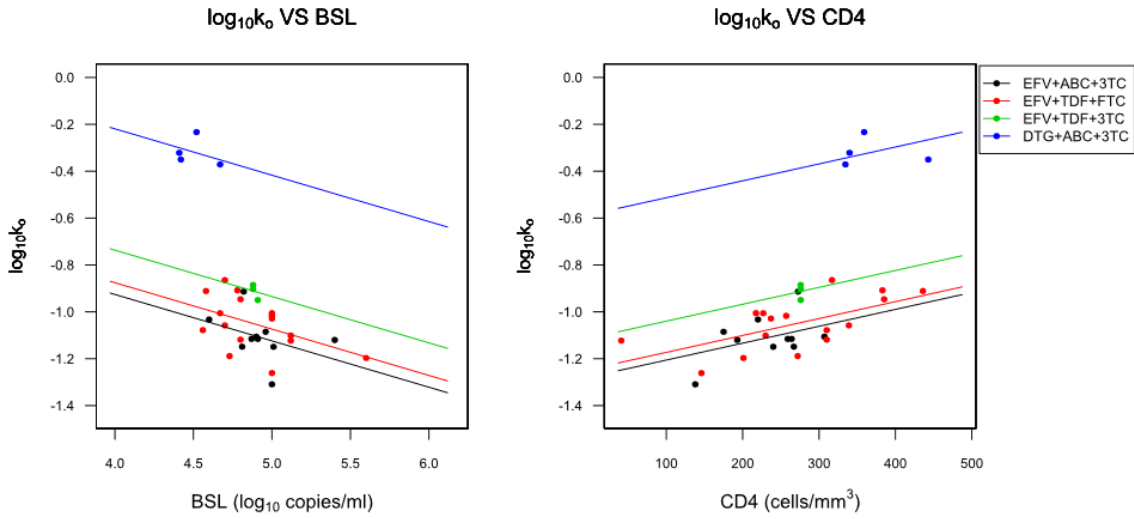
$COV_i$  : Covariate value of study arm  $i$

$COV_{med}$  : Median of the covariate



**Figure 4.** A scatter plot matrix of parameters.  $\log k_o$  and  $\log k_d$  are logarithms of the onset rate constant ( $k_o$ ) and decay rate constant ( $k_d$ ), respectively, which were estimated for the base population model by nonlinear mixed effect model analysis. BSL and CD4 are population means of mutual correlation of viral load and baseline CD4 count, respectively. When source-stratified regression analysis was conducted,  $\log k_o$  was significantly correlated to BSL and CD4 (adjusted coefficient of determination was 0.88 and 0.90, respectively). Correlation of  $\log k_d$  with BSL or CD4 was moderate (adjusted coefficient of determination was 0.67 and 0.66, respectively), although their regression coefficient was not statistically significant (the p value was 0.28 and 0.64, respectively).





**Figure 5.** Plots of regimen-strati ed regression analysis for the onset rate constant ( $\log k_o$ ).  $\log k_o$  was estimated by Bayesian post hoc analysis.

**Table 3.** Regimen-strati ed regression analysis for the onset rate constant ( $\log_{10}k_o$ )

| Regimen       | Equation                                    | Adjusted R <sup>2</sup> | Equation  | Adjusted R <sup>2</sup> |
|---------------|---|-------------------------|---|-------------------------|
| DTG+ ABC+ 3TC | $\log_{10}k_o = 0:198 \text{ BSL} + 0.572$  | 0.884                   | $\log_{10}k_o = 7:21 \cdot 10^{-4} \text{ CD4} - 0.585$ | 0.901                   |
| EFV+ ABC+ 3TC | $\log_{10}k_o = 0:198 \text{ BSL} - 0.135$  |                         | $\log_{10}k_o = 7:21 \cdot 10^{-4} \text{ CD4} - 1.28$  |                         |
| EFV+ TDF+ FTC | $\log_{10}k_o = 0:198 \text{ BSL} - 0.0845$ |                         | $\log_{10}k_o = 7:21 \cdot 10^{-4} \text{ CD4} - 1.24$  |                         |
| EFV+ TDF+ 3TC | $\log_{10}k_o = 0:198 \text{ BSL} - 0.0544$ |                         | $\log_{10}k_o = 7:21 \cdot 10^{-4} \text{ CD4} - 1.11$  |                         |

**Table 4.** Covariate screening summary

| Covariate     | parameter estimates   | OFV <sup>a</sup> | p value                    | Notes       |
|---------------|---|------------------|----------------------------|-------------|
| No covariates |   | -1158.38         |                            | Base model  |
| BSL           | $k_o = \frac{1}{2} \frac{COV_i - COV_{med}}{COV_i - COV_{med}}$ | -1170.61         | <0.05                      | Final model |
| CD4           | $k_o = \frac{1}{2} \frac{COV_i - COV_{med}}{COV_i - COV_{med}}$ | -1171.90         | <0.05                      |             |
| BSL and CD4   | $k_o = \frac{1}{2} \frac{COV_i - COV_{med}}{COV_i - COV_{med}}$ | -1175.32         | N.S. <sup>a</sup> (VS CD4) |             |

<sup>a</sup> Abbreviations: OFV, Objective function value; N.S., no signi cant.

ただし、 $COV_i$ が臨床試験の情報にない場合は、 $COV_{med}$ を $COV_i$ に代入した。この方法で、共変量モデルを作成したところ、 $k_o$ に、BSLもしくはCD4を組み込んだ場合に、有意に目的関数値 ( $OFV$ ) が減少した。BSLとCD4両方組み込んだ場合には、片方のみを組み込んだ場合と有意な差がなかったため、最終的にはより小さな $OFV$ を与えた $k_o$ にCD4を組み込んだモデルを採用した (Table 4)。解析結果は、Table 4の通りである。I-2で構築されたモデルと同様の傾向が見られ、 $\beta_1$ はドルテグラビルが最も大きく、治療効果の発現速度が最も早かった。 $\beta_2$ は1よりも大きく、CD4が少ない患者ほど治療効果は小さくなる傾向が見られた。

**Table 5.** Parameter estimates for each regimen based on a final population model

| Parameters     | Parameter estimates <sup>b</sup> |                 |                 |               |
|----------------|----------------------------------|-----------------|-----------------|---------------|
|                | DTG+ABC+3TC                      | EFV+ABC+3TC     | EFV+TDF+FTC     | EFV+TDF+3TC   |
| $E_{max}$      | 0.833 ± 0.014                    | 0.905 ± 0.052   | 0.963 ± 0.015   | 0.910 ± 0.076 |
| $\beta_1$      | 0.416 ± 0.079                    | 0.0779 ± 0.0076 | 0.0861 ± 0.0064 | 0.122 ± 0.017 |
| $\beta_2^a$    | 11.9 ± 6.3                       |                 |                 |               |
| $k_o^a$        | 0.0387 ± 0.0155                  |                 |                 |               |
| $k_d(10^{-3})$ | 0.845 ± 0.358                    | 3.93 ± 1.43     | 2.75 ± 0.48     | 1.70 ± 0.24   |
| $k_d^a$        | 0.153 ± 0.076                    |                 |                 |               |
| $t_{lag}^a$    | 1.55 ± 0.07                      |                 |                 |               |
| $\sigma_a$     | 1.96 ± 0.36                      |                 |                 |               |

<sup>a</sup>  $\beta_2$  of  $k_o$ ,  $t_{lag}$ , inter-individual variation of  $k_o$  and  $k_d$  ( $\sigma_{k_o}$  and  $\sigma_{k_d}$ ), and intra-individual variation ( $\sigma_a$ ) were assumed to be constant irrespective to any of regimens, in order to in ation of number of parameters to be estimated by curve-fitting.

<sup>b</sup> mean ± computer-calculated standard error

## I-4 考察

本研究では、MBMAを用いて、ドルテグラビルベースのレジメンとエファビレンツベースのレジメンの治療効果を定量的に比較した。MBMAは、従来のメタ解析と異なり、キネティクスに関する定量的な情報が得られることや投与量、タイムコースなどの異なるデータを統合解析できるという利点がある。特に、モデルに基づいてキネティクスに関するパラメータを得ることで、治療効果に関する多角的な理解が可能になる。また、一般的に統計技法はより多くのデータを利用することで、より確度の高い解析結果を得ることができる。したがって、MBMAは、従来のメタ解析と比較して、より信頼性の高い結果を得られることが期待される。

各レジメンの治療指数の時間推移を記述するために、単純な1次逐次反応モデルを採用した。治療指数は、被験者のうち血漿中ウイルス量が50 copies/ml以下になった患者の割合で定義しているため、 $E_{max}$ は、最大1である。解析の結果、EFV + TDF + FTCのレジメンが0.963という最も高い $E_{max}$ を示した。しかし、 $E_{max}$ の大小だけでは、レジメンの治療効果を正確に評価できない。DTG + ABC + 3TCの $E_{max}$ は0.833であり他のレジメンと比較して最も小さいが、治療効果発現速度 $k_o$ は最も大きく、治療効果減衰速度 $k_d$ は最も小さい値が得られている。したがって、レジメンの評価には、これら三つのパラメータをそれぞれ評価し、その治療効果の特性を理解することが必要である。なお、臨床試験では、たとえ治療効果が認められても、患者個人の理由や副作用の発現によって離脱する患者が存在する。 $k_d$ は、単純な治療効果の減衰のみではなく、こういった臨床試験の途中離脱も含意したパラメータである点には解釈上の注意が必要である。

エファビレンツベースのレジメンの三つを比較すると、EFV + ABC + 3TCのレジメンが最も治療効果が小さかった( $E_{max}$ 最小、 $k_o$ 最小、 $k_d$ 最大)。実際、EFV + ABC + 3TCのレジメンは、EFV + TDF + FTCと比較して、副作用が大きいという報告がある[16, 19]。また、EFV + TDF + FTC、EFV + TDF + 3TCのレジメンは、HBV感染、結核感染、妊娠女性での利用がWHOで推奨されている[6]。こういった背景を踏まえると、EFV + ABC + 3TCのレジメンが、他のレジメンと比較すると相対的に安全性が低いことが推察され、それが本モデルでの治療効果の減衰速度を表す $k_d$ が大きいことに反映されている可能性がある。

DTG + ABC + 3TCは、エファビレンツベースのレジメンと比較して、 $k_o$ が4倍以上大きく、 $k_d$ が約1/2以下になった。このことは、治療発現速度が速く、治療減衰速度が遅いことを示唆している。治療発現速度 $k_o$ が大きいことは、細胞系での実験結果とも対応している[46]。 $k_d$ が小さいことに関しては、薬剤の薬剤耐性ウイルスに対する頑

健性と副作用の発現率という二つの観点から説明可能だと考えられる。一つ目の説明として、一般に治療においては、薬剤耐性ウイルスの出現が問題となるが、ドルテグラビルの場合は、genetic barrier(薬剤耐性を獲得するのに必要なアミノ酸変異の数)が高く、交叉耐性の出現が少ないとの報告がある[47, 48]。その一方で、エファビレンツの場合には、genetic barrierが低く、交叉耐性の出現が多いとの報告がある[48, 49]。したがって、こういった薬剤の特性から、薬剤耐性ウイルスの出現による治療効果減退のリスクは両薬剤とも抱えているものの、ドルテグラビルの方が相対的に長く治療効果が持続させられる可能性がある。また、もう一つの理由として、前段落のエファビレンツ同士のレジメン比較の議論と同様に、副作用発現による離脱が少ないことを反映している可能性も考えられる。実際に、エファビレンツベースのレジメンと比較して、ドルテグラビルベースのレジメンの方が副作用が少ないという報告もされている[50]。これら二つの理由から、本研究で推定された $k_d$ がドルテグラビルの方が小さいことは、合理的な結果であったと考えられる。

共変量モデルの解析の結果から、レジメンによらず治療前の正常CD4のカウント数と $k_0$ との正の相関がみられた。この結果は、正常CD4のカウント数が少ないほど治療効果の発現が遅いというSkowronら[51]の報告と一致している。ここで、HIV抗ウイルス薬の評価では、治療効果をアーム内で治療効果のあった患者数の割合で定義する、すなわち集団としての治療効果を評価するため、共変量であるCD4カウントに関してもアーム内での中央値あるいは平均値でしか解析に利用できない。したがって、本研究の結果が個人レベルの解析においても同じ傾向が見られるかは未だ不明である。

以上のように、MBMAを用いて複数のデータソースを統合し解析した結果、ドルテグラビルベースのレジメンが、エファビレンツベースのレジメンと比較して、治療発現速度が速く持続性が高いことが示唆された。

## 第二章

### HIVの薬剤耐性に関する予測モデル構築法の提案と検証

HIV治療において、薬剤耐性ウイルスは、治療予後に大きな影響を及ぼす[52, 53]。ウイルスが耐性を示す薬剤を遺伝子型から予測し、適切な治療薬を選択することは、治療の効率化に不可欠である。こういった予測モデルの構築に関する研究が活発に進められており[54{58]、その背景として、薬剤耐性ウイルスのタンパク質情報と薬剤耐性の実験値を収集したStanford HIV Drug resistance databaseに代表されるようなデータベースが成熟してきた経緯がある[55, 59]。

これらの予測モデルは、タンパク質の一次配列をone-hot encoding等で数値配列に変換し、その数値配列を説明変数、薬剤耐性を目的変数として、各種機械学習法(サポートベクターマシン、ディープラーニングなど)により構築された分類モデルあるいは回帰モデルである。Margaretら[54]は、近年、幅広い分野で成果を上げているDeep Learningを利用して、分類問題に取り組み、正解率(Accuracy)は0.9程度と高い精度を出すことに成功し、薬剤耐性に寄与するアミノ酸残基も同定している。また、Nikoら[55]が開発したgeno2phenoでは、サポートベクターマシンを利用して、決定係数0.7程度の回帰モデルが得られている。これらのモデルは、インターネット上で公開され、臨床の現場で利用されている。

しかしながら、これらのモデルには、次のような問題がある。一つは、タンパク質の一次配列に基づく予測であるため、構造バイオインフォマティクスで重視されるタンパク質の立体構造情報が欠落している点である[60]。タンパク質の機能は三次構造で決定されるのは薬剤耐性の獲得でも同様である。二つは、配列未確定データの処理に関する問題である。配列未確定データとは、一つのアミノ酸残基の位置に複数候補のアミノ酸が存在する配列のことである。これらのデータを削除するか、すべてのパターンを考え利用するかは著者により見解が異なる。前者の場合[54]には、サンプルサイズの問題により学習が不十分となる可能性があり、後者[57, 58]の場合には、不確定なデータから生成した膨大なデータに対してモデルが過度に適合してしまう可能性がある。

そこで、本研究では、これらの問題を解決するために、ホモロジーモデリングとComparative Molecular Field Analysis (CoMFA)の手法を応用して、タンパク質の一次配列情報をタンパク質の三次元構造情報へと変換した。CoMFAは、Cramerら[61]によって開発された3D-QSARの手法であり、目的のタンパク質を格子で囲い込み、各格子点のプローブ原子とタンパク質との相互作用(立体・静電ポテンシャル)を計算する。通常低分子薬物の3次元構造活性相関に利用される方法であるが、これを利用してタンパ

ク質3次元構造を物理化学的に表現する数値配列情報に変換した。一方、配列未確定データでは、すべての組み合わせを数え上げた後、それぞれの存在確率を推定することによって、データに重みづけを行った。機械学習法を予測モデルの構築に利用することとし、Partial Least Squares (PLS)[62], Random Forest (RF)[63], LightGBM[64], Support Vector regression (SVR)[65]の各手法を採用し、それぞれの精度を評価した。さらに、薬剤耐性に寄与するタンパク質構造に関連する重要なパラメータを空間上で可視化した。これらの解析法を比較的タンパクサイズが小さく、計算コストの小さいプロテアーゼ阻害剤に対して適用した。

## II - 1 データ収集とデータの前処理

Stanford HIV Drug Resistance database[59, 66] から、プロテアーゼ阻害剤の薬剤耐性インデックスを取得した (Table 6)。本データベース上の薬剤耐性インデックスとは、変異体の HIV の IC50 を野生型の HIV の IC50 で除した Fold change(FC) である。一般に、 $FC \geq 3.5$  を薬剤耐性あり、 $FC < 3.5$  を薬剤耐性なしとみなされている [54, 67]。データベースには、タンパク質の一次配列が完全に決定しているデータ (配列確定データ) とタンパク質の一次配列が未確定部分のあるデータ (以下では、配列未確定データとする) の二種類が混在している。配列未確定データとは、配列中のあるアミノ酸残基の位置に複数候補のアミノ酸が存在するデータのことである。

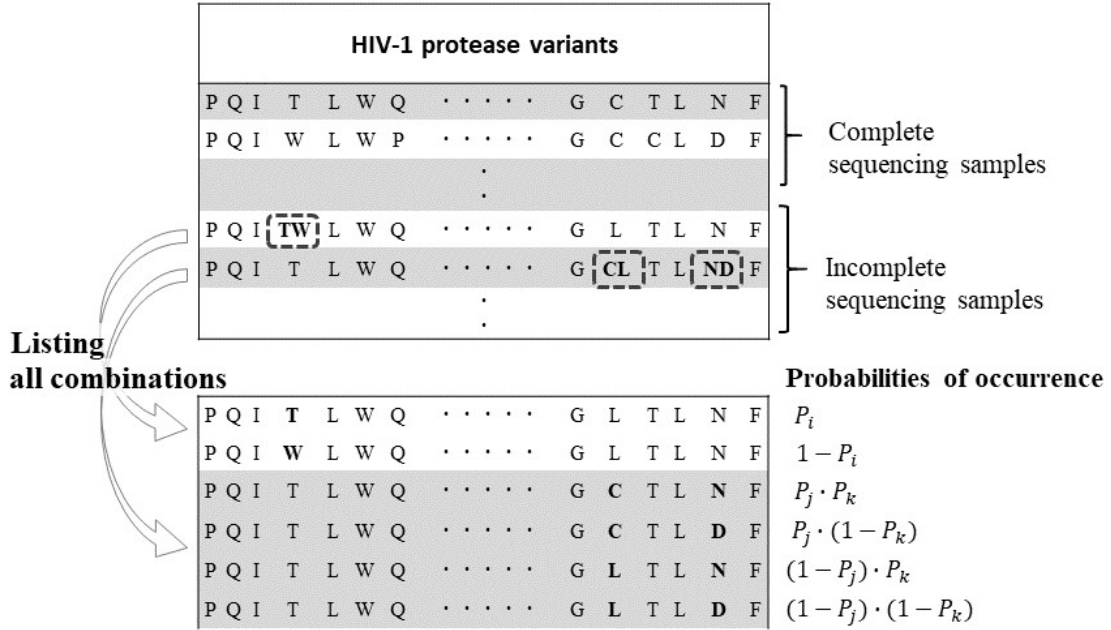
**Table 6.** Numbers of sequencing samples listed in the database for each HIV-1 protease inhibitors.

| Drug          | Number of complete sequencing samples | Incomplete sequencing samples |  |  |
|---------------|---------------------------------------|-------------------------------|--|--|
|               |                                       | Number of samples             | Total sum of all combinations <sup>a</sup> | Practical total sum of all combinations <sup>b</sup> |
| Atazanavir    | 463                                   | 596                           | 219847                                     | 407.97   |
| Darunavir     | 264                                   | 401                           | 210578                                     | 275.5  |
| Fosamprenavir | 726                                   | 834                           | 232626                                     | 581.63   |
| Indinavir     | 759                                   | 849                           | 233116                                     | 597.86   |
| Lopinavir     | 600                                   | 773                           | 231749                                     | 536.13   |
| Nel navir     | 781                                   | 874                           | 234104                                     | 613.68   |
| Saquinavir    | 758                                   | 846                           | 233304                                     | 592.94   |
| Tipranavir    | 302                                   | 464                           | 212140                                     | 316.77   |

<sup>a</sup> The total sum of all possible combinations of sequences in each incomplete sequencing sample.

<sup>b</sup> The total sum of all possible combinations weighted by their probability of occurrence in each incomplete sequencing sample.

配列未確定データでは、アミノ酸残基に複数候補が考えられるため、全ての組み合わせを調べ上げ、それらの存在確率を推定する必要がある (Figure 6)。



**Figure 6.** Scheme of treatment of incomplete sequencing samples. All possible combinations were listed from incomplete sequencing samples and their information was weighted by the conditional probability of occurrence. The conditional probability of occurrence of each amino acid at each position was determined based on amino acid sequences of complete sequencing samples.

そこで、未確定の位置*i*における各アミノ酸候補の出現確率を、異なる*m*個の位置のアミノ酸の種類に基づいて計算される条件付確率とした。これは、ウイルス遺伝子の変異が集団の遺伝的特性の変化として捉えられる集団遺伝学的発想に基づいている [68]。アミノ酸の種類を確率変数として考え、アミノ酸配列の位置*i*のアミノ酸を  $X_i \in \{A, R, N, \dots, Y\}$  とし、周辺確率分布関数  $P(X_i)$ 、同時確率分布関数  $P(X_i; X_j)$  とするとき、情報エントロピー、相互情報量、標準化相互情報量は次のようにあらわされる [69]。

$$H(X_i) = - \sum_{x_i \in X_i} p(x_i) \log p(x_i) \quad (6)$$

$$I(X_i; X_j) = - \sum_{x_i \in X_i} \sum_{x_j \in X_j} p(x_i; x_j) \log \frac{p(x_i; x_j)}{p(x_i) p(x_j)} \quad (7)$$

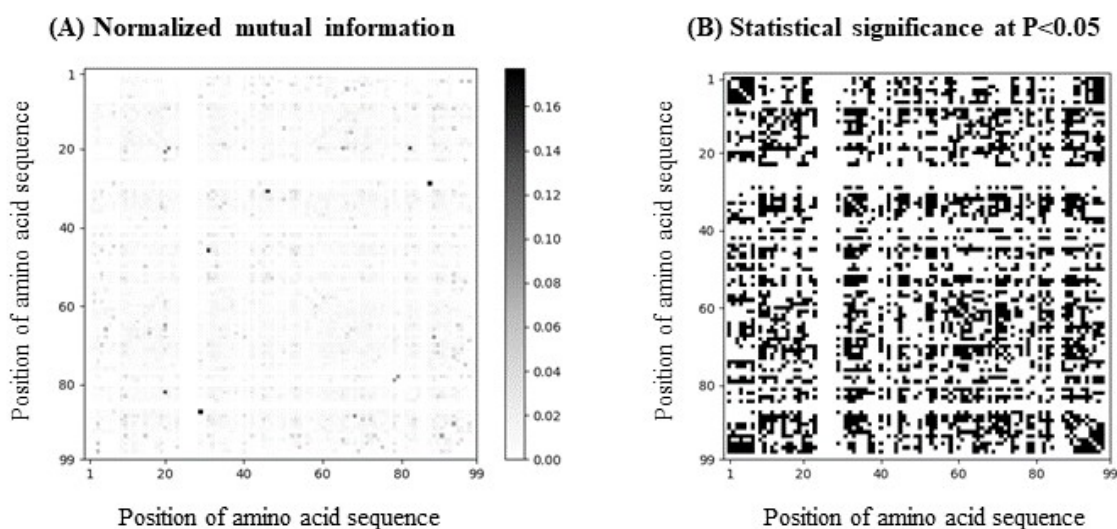
$$NMI(X_i; X_j) = \frac{I(X_i; X_j)}{H(X_i) H(X_j)} \quad (0 \leq NMI(X_i; X_j) \leq 1) \quad (8)$$

まず、配列確定データを利用して、NMI行列を計算した。次に、各行列要素に対して permutation test (P値 0.05) [70] により有意性を確認した。Figure 7は、NMI行列の濃淡表



示および有意な相関を示す対となる位置を示したものである。次に、各アミノ酸変異位置に対して相関性のある位置の組み合わせ  $\mathbf{X}_m$  を決定した。ただし、 $\mathbf{X}_m$  の最大長は10とし、それを超える場合はNMIが大きい上位10個の位置を選択した。

得られたNMIから、配列未確定データにおける全組み合わせ候補それぞれの存在確率を推定した。Table 6には、配列未確定データから書き出したすべての組み合わせの総数およびそれぞれに存在確率を乗じた期待度数の総和（実質的データ量）を記載した。実質的な情報量は配列未確定データ数よりも小さくなったが、これは配列確定データ中に参照できる配列が存在せず、結果として存在確率が0となった配列も含まれるためである。最終的に、配列確定データのみの場合と比較して、約2倍程度のデータ量が確保できた。



**Figure 7.** Normalized mutual information on co-occurrence of amino acids between any two positions. (A) Gray-scale image matrix of normalized mutual information (NMI). (b) Statistical significance in NMI determined by permutation test at  $P < 0.05$  (black).

## II - 2 分子場ポテンシャルの計算

ホモロジーモデリング法[71]により各プロテアーゼ変異体の立体構造を一次配列に基づいてシミュレートした。各プロテアーゼ変異体に対して、ホモロジーモデリングを行う際のテンプレートには、立体構造既知のプロテアーゼ変異体のうち一次配列の類似度が最も高いものを採用した。ホモロジーモデリングにより得られたプロテアーゼ変異体の立体構造は、最もよくテンプレートとして使われたものに重ね合わせた。テンプレートとして利用したプロテアーゼの立体構造は、Table 7の通りである[72][83]。

推定された立体構造を格子点間隔が2Åで十分な大きさの格子内に置き、各格子点でのポテンシャルエネルギーを計算した。エネルギー計算には、sp<sup>3</sup>炭素と同じファンデルワールス半径を持ち、原子電荷が+1のプローブ原子を仮定した。その上で、プローブ原子とプロテアーゼとの間に働く立体相互作用および静電相互作用としてそれぞれファンデルワールスポテンシャル[84]、クーロンポテンシャル[85]を計算し、これらの分子場ポテンシャルを各プロテアーゼ変異体の立体構造を表す記述子(説明変数)として利用した。

**Table 7.** Template proteins for homology modeling of HIV protease variants

| Drug          | Template (PDB ID) |      |      |      |
|---------------|-------------------|------|------|------|
| Atazanavir    | <b>3EKW</b>       | 3KEY |      |      |
| Darunavir     | <b>3JVY</b>       | 3JW2 |      |      |
| Fosamprenavir | <b>3NU9</b>       | 3NUJ | 3NUO |      |
| Indinavir     | <b>1SDT</b>       | 1SDV | 1SGU | 1HSG |
| Lopinavir     | <b>6DJ1</b>       | 6DJ2 | 1MUI |      |
| Nel navir     | <b>2PYM</b>       | 2PYN | 2Q63 | 2Q64 |
| Saquinavir    | <b>3D1Y</b>       |      |      |      |
| Tipranavir    | <b>2O4P</b>       | 2O4N |      |      |

There are several available complex structures with HIV protease variants for each drug. For each protease variant to be subjected to homology modeling, one of the protein structures listed was selected as a template according to the similarity of primary amino acid sequences.

## II - 3 予測モデルの構築

全体データの80%を学習用、20%を評価用(External)にランダムに分割した。以下の特徴抽出およびハイパーパラメータの選択は学習用データの80%を使って行うこととし、最終的な予測モデルの評価は残りの20%のデータを用いて行った。モデル開発では、分

子場ポテンシャルを予測子として薬剤耐性インデックスを定量的に予測する回帰モデルを基本とし、Partial Least Squares (PLS), Random Forest (RF), LightGBM (LGBM), Support Vector regression (SVR) といった各種機械学習アルゴリズムを利用した。

回帰モデルの損失関数 (Loss function)、評価関数 (Evaluation function) を以下のように定義した。

$$Lossfunction = \sum_P \sum_P p_{ij} (y_{obs:i} - y_{pred:ij})^2 \quad (9)$$

$$Evaluationfunction(R^2) = 1 - \frac{\sum_P p_{ij} (y_{obs:i} - y_{pred:ij})^2}{(\sum_P y_{obs:i} - \overline{y_{obs}})^2} \quad (10)$$

$y_{obs:i}$  and  $\overline{y_{obs}}$  were an observed log  $FC$  of  $IC_{50}$  for sample  $i$  and their average, respectively; and  $y_{pred:ij}$  and  $p_{ij}$  were predicted log  $FC$  and conditional probability for the  $j$ -th candidate of sample  $i$ , respectively.

## II - 3 - a 特徴抽出

特徴抽出には配列確定データのみを利用した。分子場ポテンシャルは、プローブ原子と標的タンパク質上の原子の中心距離が近接すると発散する。そこで、各分子場ポテンシャルに対して、全データの5パーセンタイル、95パーセンタイルをそれぞれ下限値、上限値とし、それ以下あるいはそれ以上は限度値で打ち切った。また、プロテアーゼ変異体内での密度分布が正規分布から極端に外れる分子場ポテンシャル、具体的には歪度の絶対値が2.5以上の分子場ポテンシャルは予測子候補から取り除いた。さらに、変異体内での標準偏差2 kcal/mol以下で変動が少ない分子場ポテンシャルも予測子候補から取り除いた [86]。

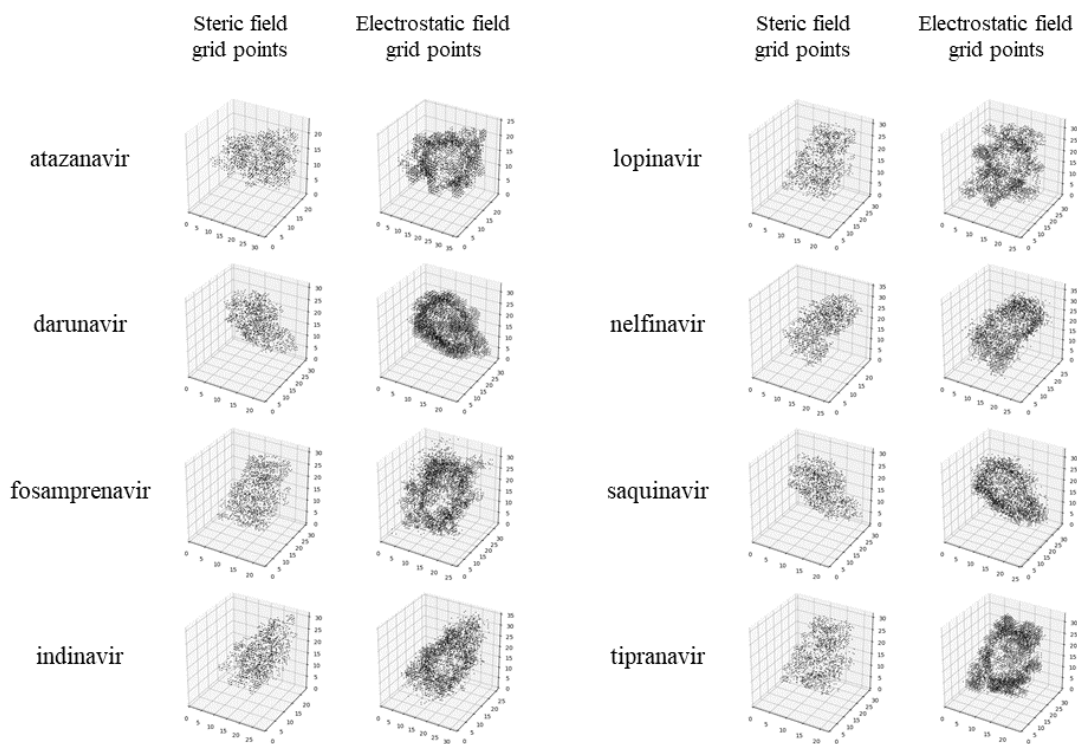
次に、線形SVRカーネルを利用して分子場ポテンシャルの重要度を推定し、Recursive feature elimination (RFE)[87]により説明力のない分子場ポテンシャルを削減した。RFEは分子場ポテンシャルであるファンデルワールスポテンシャルとクーロンポテンシャルのそれぞれに独立して適用し、ファンデルワールスポテンシャルでは、重要度の低い記述子から順に2個ずつ減らし、最終的に1/2になるまで繰り返したのに対し、クーロンポテンシャルでは、10個ずつ記述子を削減して1/8にした。結果として、各薬剤に残された特徴量の数は、約4000であった (Table 8)。Figure 8は、削除されずに残った有効な分子場ポテンシャルの座標を空間上に示したものであり、プロテアーゼタンパク質の構造を取り囲むような点が抽出されていることが理解された。

**Table 8.** Number of molecular field energies selected as a feature of machine learning

| Drug          | Grid size <sup>a</sup><br>(A) | Before feature selection |                            |        | After feature selection |                            |       |
|---------------|-------------------------------|--------------------------|----------------------------|--------|-------------------------|----------------------------|-------|
|               |                               | Steric <sup>b</sup>      | Electrostatic <sup>b</sup> | Total  | Steric <sup>b</sup>     | Electrostatic <sup>b</sup> | Total |
| Atazanavir    | 39 × 27 × 28                  | 29,484                   | 29,484                     | 58,968 | 1,223                   | 2,728                      | 3,951 |
| Darunavir     | 27 × 35 × 36                  | 34,020                   | 34,020                     | 68,040 | 1,094                   | 3,443                      | 4,537 |
| Fosamprenavir | 27 × 34 × 35                  | 32,130                   | 32,130                     | 64,260 | 1,210                   | 3,268                      | 4,478 |
| Indinavir     | 34 × 27 × 36                  | 33,048                   | 33,048                     | 66,096 | 1,121                   | 2,549                      | 3,670 |
| Lopinavir     | 27 × 34 × 36                  | 33,048                   | 33,048                     | 66,096 | 1,096                   | 2,664                      | 3,760 |
| Nelfinavir    | 30 × 28 × 40                  | 33,600                   | 33,600                     | 67,200 | 1,199                   | 2,344                      | 3,543 |
| Saquinavir    | 27 × 35 × 35                  | 33,075                   | 33,075                     | 66,150 | 1,107                   | 2,362                      | 3,469 |
| Tipranavir    | 26 × 35 × 36                  | 32,760                   | 32,760                     | 65,520 | 1,127                   | 3,285                      | 4,412 |

<sup>a</sup> Grid sizes were set to embed an overlaid aggregate of all HIV protease variants.

<sup>b</sup> Numbers of steric and electrostatic potential energies.

**Figure 8.** Selected grid points of steric (left) and electrostatic (right) molecular field parameters in the analysis of drug resistance for each HIV-1 protease inhibitors. The grid points were selected by preprocessing and SVR feature selection.

## II - 3 - b ハイパーパラメータの最適化

各種機械学習には、ハイパーパラメータと呼ばれる学習プロセスを制御するパラメータが存在する。この設定によってはモデルの精度やパフォーマンスが大きく左右されることがある。そこで、PLSを除く各機械学習(SVR、RF、LGBM)では、Tree-structured Parzen Estimator Approach(TPE)[88]を利用して、ハイパーパラメータの最適化を行った。最適化するハイパーパラメータの候補は、Table 9の通りである。データは三分割し、一つをバリデーション用のデータセット(バリデーションセット)とし、残りの二つをトレーニング用のデータセット(トレーニングセット)として、3-fold cross validation[89]を行った。トレーニングセットのうち、II-3-a特徴抽出の項で抽出された分子場ポテンシャルのみを予測子として利用し、学習の前にそれらを正規化した。バリデーションセットの正規化にも、トレーニングセットで定めた平均、分散を利用した。各機械学習法による学習過程では、トレーニングセットの損失関数(式(9))を最小化するように学習が進み、最後にバリデーションセットで評価関数(式(10))を計算した。クロスバリデーションを3回繰り返し、評価関数の三回の平均値をハイパーパラメータの評価値とした。その評価値に基づいて、TPEのアルゴリズムに従ってハイパーパラメータを30回更新した。一方、PLSのハイパーパラメータである主成分数の決定では、主成分を順次増やしながらかrossvalidationによる評価値(式(10))を計算し、評価値が最大となる最も精度の良い主成分数を取得した。Table 10は、最終的に得られたハイパーパラメータ、PLSでの次元数をまとめたものである。

**Table 9.** Ranges of hyperparameters in model optimization

| LightGBM         |   |    | Random Forest Regression |               |     | Support Vector Regression |               |   |
|------------------|---|----|--------------------------|---------------|-----|---------------------------|---------------|---|
| num_leaves       | 2 | 60 | Bootstrap                | true or false |     | kernel                    | linear or rbf |   |
| n_estimators     | 2 | 60 | max_depth                | 5             | 150 | C                         | 0             | 4 |
| bagging_fraction | 0 | 1  | max_features             | auto or sqrt  |     | epsilon                   | 0             | 1 |
|                  |   |    | min_samples_leaf         | 1             | 10  |                           |               |   |
|                  |   |    | min_samples_split        | 2             | 10  |                           |               |   |
|                  |   |    | n_estimators             | 5             | 150 |                           |               |   |

LightGBM model were constructed using LightGBM 2.3.0, while random forest regression and support vector regression using scikit-learn 0.23.1.

**Table 10.** Optimized hyperparameters of LightGBM (LGBM), Random Forest Regression (RF), Support Vector Regression (SVR), and Partial Least Squares (PLS) models for each drug

| Model | Hyperparameter                 | Drug <sup>a</sup> |        |        |        |        |        |        |        |
|-------|--------------------------------|-------------------|--------|--------|--------|--------|--------|--------|--------|
|       |                                | ATV               | DRV    | FPV    | IDV    | LPV    | NFV    | SQV    | TPV    |
| LGBM  | num_leaves                     | 9                 | 12     | 12     | 13     | 9      | 17     | 8      | 20     |
|       | n_estimators                   | 56                | 60     | 58     | 60     | 60     | 60     | 53     | 52     |
|       | bagging_fraction               | 0.19              | 0.482  | 0.012  | 0.993  | 0.98   | 0.293  | 0.779  | 0.996  |
| RF    | bootstrap                      | TRUE              | TRUE   | TRUE   | TRUE   | TRUE   | TRUE   | TRUE   | TRUE   |
|       | max_depth                      | 864               | 345    | 998    | 178    | 756    | 386    | 453    | 624    |
|       | max_features                   | auto              | auto   | auto   | auto   | auto   | auto   | auto   | auto   |
|       | min_samples_leaf               | 7                 | 2      | 2      | 2      | 8      | 3      | 3      | 6      |
|       | min_samples_split              | 10                | 10     | 10     | 10     | 10     | 10     | 7      | 10     |
|       | n_estimators                   | 102               | 108    | 150    | 50     | 105    | 12     | 11     | 6      |
| SVR   | kernel                         | linear            | linear | linear | linear | linear | linear | linear | linear |
|       | C                              | 0.002             | 3.053  | 1.075  | 3.461  | 0.002  | 3.361  | 0.012  | 0.861  |
|       | epsilon                        | 0.328             | 0.256  | 0.477  | 0.004  | 0.332  | 0.467  | 0.348  | 0.34   |
| PLS   | number of principal components | 41                | 111    | 34     | 60     | 98     | 113    | 28     | 112    |

<sup>a</sup> Abbreviations: ATV, atazanavir; DRV, darunavir; FPV, fosamprenavir; IDV, indinavir; LPV, lopinavir; NFV, nel navir; SQV, saquinavir; TPV, tipranavir.

### II - 3 - c モデルの最終評価

前述の通り、全体データの80%をトレーニング用データ(トレーニングセット)に、20%を外部評価用データ(テストセット)とした。トレーニングセットでは、II - 3 - aで抽出された分子場ポテンシャルを正規化し、その正規化に利用したパラメータでテストセットも同様に正規化した。各種機械学習法(SVR、RF、LGBM、PLS)に対して、II - 3 - bのクロスバリデーションで決定されたハイパーパラメータを適用して、トレーニングセット全体を再学習した(Table 11)。トレーニングによって得られた最終予測モデルの予測精度をテストセットで評価した(Table 12)。精度評価には、重み付き決定係数(式(10))を利用しており、1に近いほど高い予測精度と解釈される。Table 12のテスト結果での精度は、チプラナビルを除いて良好であった。LightGBMを用いた場合に最も高い予測精度が得られる傾向にあったが、機械学習のモデルごとにそれほど大きな性能の差は観察されなかった。

予測モデルの構築においては、薬剤耐性があるかないかを識別するための分類モデルも広く提案されている。そこで、本研究でも回帰モデルの精度に加えて、既報の手法との比較をするために分類性能を評価することとした。既報と同様に、FC=3.5を境界値として、 $FC < 3.5$ を薬剤耐性無し、 $FC \geq 3.5$ を薬剤耐性ありとしてデータを扱い、先のテストセットで分類性能を評価した(Table 13)。分類精度の評価には、既報の論文でも使用されていたaccuracy、precision、true-positive rate (TPR)、true-negative rate (TNR)、false-positive rate (FPR)、false-negative rate (FNR)、area under the ROC curve (AUC)、F1 scoreの複数のメトリクスを利用した[90, 91]。回帰モデルでの精度と同様に、予測モデルで最も予測精度が高かったのは、LightGBMを用いた場合であったものの、機械学習法ごとの性能差はほとんど認められなかった。

**Table 11.** Weighted determination coefficients for prediction in training dataset ( $R^2$ )

| Drug          | Weighted determination coefficient |                          |                           |                       |
|---------------|------------------------------------|--------------------------|---------------------------|-----------------------|
|               | LightGBM                           | Random Forest Regression | Support Vector Regression | Partial Least Squares |
| Atazanavir    | 0.948                              | 0.954                    | 0.888                     | 0.853                 |
| Darunavir     | 0.981                              | 0.951                    | 0.911                     | 0.901                 |
| Fosamprenavir | 0.929                              | 0.946                    | 0.822                     | 0.778                 |
| Indinavir     | 0.951                              | 0.958                    | 0.986                     | 0.839                 |
| Lopinavir     | 0.952                              | 0.968                    | 0.921                     | 0.896                 |
| Nel navir     | 0.957                              | 0.925                    | 0.831                     | 0.819                 |
| Saquinavir    | 0.905                              | 0.945                    | 0.892                     | 0.748                 |
| Tipranavir    | 0.975                              | 0.827                    | 0.781                     | 0.803                 |

**Table 12.** Weighted determination coefficients for prediction in external dataset ( $R^2$ )

| Drug          | Weighted determination coefficient |                          |                           |                       |
|---------------|------------------------------------|--------------------------|---------------------------|-----------------------|
|               | LightGBM                           | Random Forest Regression | Support Vector Regression | Partial Least Squares |
| Atazanavir    | 0.792                              | 0.770                    | 0.737                     | 0.759                 |
| Darunavir     | 0.749                              | 0.694                    | 0.711                     | 0.693                 |
| Fosamprenavir | 0.718                              | 0.667                    | 0.683                     | 0.701                 |
| Indinavir     | 0.830                              | 0.790                    | 0.783                     | 0.806                 |
| Lopinavir     | 0.862                              | 0.810                    | 0.829                     | 0.837                 |
| Nel navir     | 0.760                              | 0.743                    | 0.714                     | 0.705                 |
| Saquinavir    | 0.776                              | 0.632                    | 0.748                     | 0.676                 |
| Tipranavir    | 0.513                              | 0.491                    | 0.470                     | 0.488                 |



**Table 13.** Goodness of classification of LightGBM, Random Forest Regression, Support Vector Regression, and Partial Least Squares models<sup>a</sup>

|                                  | Drug <sup>b</sup> | Accuracy | Precision | TPR <sup>c</sup> | TNR <sup>c</sup> | FPR <sup>c</sup> | FNR <sup>c</sup> | AUC <sup>c</sup> | F1 score |
|----------------------------------|-------------------|----------|-----------|------------------|------------------|------------------|------------------|------------------|----------|
| <b>LightGBM</b>                  | ATV               | 0.914    | 0.881     | 0.947            | 0.883            | 0.117            | 0.0533           | 0.915            | 0.913    |
|                                  | DRV               | 0.921    | 0.810     | 0.805            | 0.951            | 0.0487           | 0.195            | 0.878            | 0.807    |
|                                  | FPV               | 0.896    | 0.789     | 0.921            | 0.885            | 0.115            | 0.0795           | 0.903            | 0.850    |
|                                  | IDV               | 0.924    | 0.897     | 0.917            | 0.928            | 0.0716           | 0.0834           | 0.922            | 0.907    |
|                                  | LPV               | 0.915    | 0.895     | 0.929            | 0.903            | 0.0972           | 0.0711           | 0.916            | 0.912    |
|                                  | NFV               | 0.906    | 0.888     | 0.938            | 0.873            | 0.127            | 0.0624           | 0.905            | 0.912    |
|                                  | SQV               | 0.909    | 0.832     | 0.929            | 0.898            | 0.102            | 0.0712           | 0.913            | 0.878    |
|                                  | TPV               | 0.890    | 0.681     | 0.567            | 0.950            | 0.0497           | 0.433            | 0.759            | 0.619    |
| <b>Random Forest Regression</b>  | ATV               | 0.894    | 0.840     | 0.961            | 0.833            | 0.167            | 0.0388           | 0.897            | 0.897    |
|                                  | DRV               | 0.921    | 0.889     | 0.701            | 0.977            | 0.0227           | 0.299            | 0.839            | 0.784    |
|                                  | FPV               | 0.869    | 0.744     | 0.901            | 0.855            | 0.145            | 0.0992           | 0.878            | 0.815    |
|                                  | IDV               | 0.905    | 0.847     | 0.934            | 0.885            | 0.115            | 0.0663           | 0.909            | 0.888    |
|                                  | LPV               | 0.910    | 0.878     | 0.940            | 0.884            | 0.116            | 0.0599           | 0.912            | 0.908    |
|                                  | NFV               | 0.899    | 0.873     | 0.941            | 0.854            | 0.146            | 0.0585           | 0.898            | 0.906    |
|                                  | SQV               | 0.904    | 0.807     | 0.957            | 0.875            | 0.125            | 0.0433           | 0.916            | 0.876    |
|                                  | TPV               | 0.868    | 0.624     | 0.405            | 0.954            | 0.0457           | 0.595            | 0.680            | 0.491    |
| <b>Support Vector Regression</b> | ATV               | 0.924    | 0.881     | 0.973            | 0.880            | 0.120            | 0.0275           | 0.926            | 0.924    |
|                                  | DRV               | 0.954    | 0.900     | 0.874            | 0.975            | 0.0250           | 0.126            | 0.925            | 0.887    |
|                                  | FPV               | 0.879    | 0.800     | 0.830            | 0.903            | 0.0975           | 0.170            | 0.866            | 0.814    |
|                                  | IDV               | 0.935    | 0.906     | 0.937            | 0.934            | 0.066            | 0.0630           | 0.935            | 0.921    |
|                                  | LPV               | 0.929    | 0.894     | 0.964            | 0.898            | 0.102            | 0.0362           | 0.931            | 0.928    |
|                                  | NFV               | 0.87     | 0.850     | 0.909            | 0.828            | 0.172            | 0.0913           | 0.868            | 0.878    |
|                                  | SQV               | 0.894    | 0.802     | 0.93             | 0.875            | 0.125            | 0.0696           | 0.903            | 0.861    |
|                                  | TPV               | 0.878    | 0.588     | 0.755            | 0.901            | 0.0992           | 0.245            | 0.828            | 0.661    |
| <b>Partial Least Squares</b>     | ATV               | 0.918    | 0.863     | 0.985            | 0.857            | 0.143            | 0.0154           | 0.921            | 0.920    |
|                                  | DRV               | 0.897    | 0.786     | 0.683            | 0.952            | 0.0480           | 0.317            | 0.818            | 0.731    |
|                                  | FPV               | 0.862    | 0.777     | 0.797            | 0.893            | 0.107            | 0.203            | 0.845            | 0.787    |
|                                  | IDV               | 0.911    | 0.880     | 0.904            | 0.916            | 0.0838           | 0.096            | 0.910            | 0.892    |
|                                  | LPV               | 0.914    | 0.888     | 0.936            | 0.894            | 0.106            | 0.064            | 0.915            | 0.911    |
|                                  | NFV               | 0.897    | 0.871     | 0.941            | 0.850            | 0.150            | 0.059            | 0.896            | 0.904    |
|                                  | SQV               | 0.871    | 0.780     | 0.882            | 0.865            | 0.135            | 0.118            | 0.873            | 0.828    |
|                                  | TPV               | 0.884    | 0.612     | 0.724            | 0.914            | 0.0860           | 0.276            | 0.819            | 0.663    |

<sup>a</sup>Goodness of classification was evaluated upon defining the threshold as a fold change of 3.5.

<sup>b</sup>Abbreviations: ATV, atazanavir; DRV, darunavir; FPV, fosamprenavir; IDV, indinavir; LPV, lopinavir; NFV, nel navir; SQV, saquinavir; TPV, tipranavir.

<sup>c</sup>Abbreviations: TPR, true positive ratio; TNR, true negative ratio; FPR, false positive ratio; FNR, false negative ratio; AUC, area under the ROC curve.

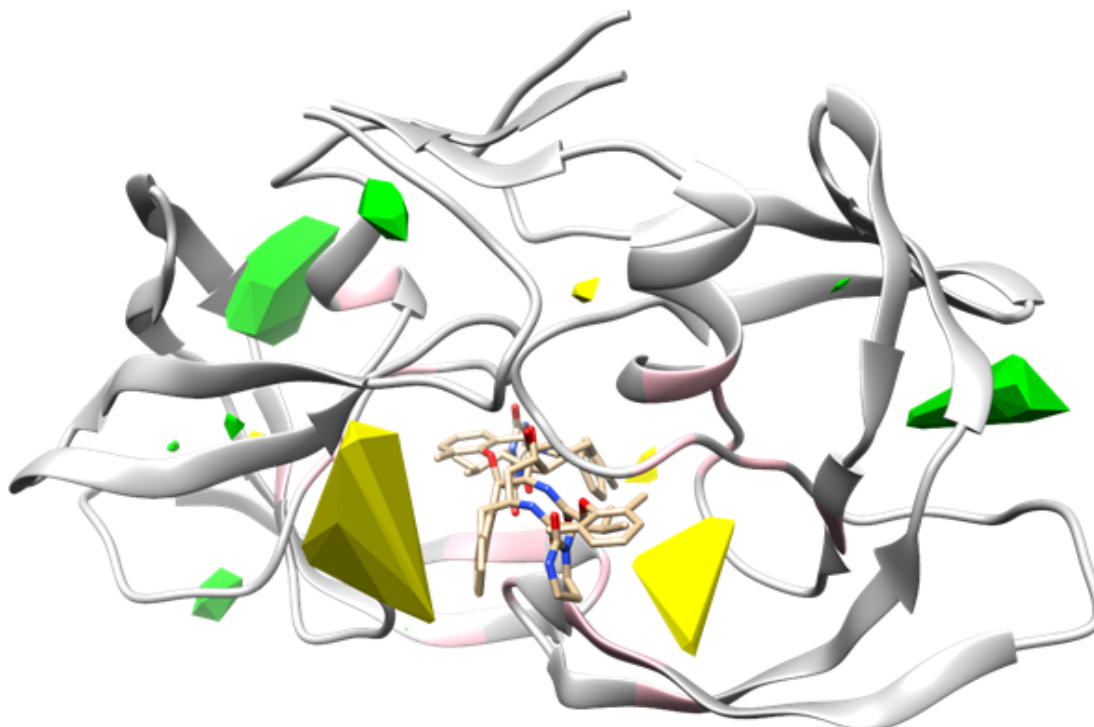
## II - 4 薬剤耐性ウイルスの構造要因解析

CoMFA[61]は、一般に、基本骨格を同一とする低分子化合物の構造活性相関解析に用いられる3D-QSARの手法である。これに対し、本モデルは、変異したタンパク質に対する薬剤の効果(薬剤耐性)をシミュレートするものである。これは、リガンド構造に基づく物理化学的性質の違いで活性が変化するというCoMFAの概念を逆転させ、レセプター側(タンパク質)の物性変化が活性を変化させるという発想に基づいている。なお、CoMFAには、予測対象である生物活性に対して重要となる変動要因を可視化できるという特長がある。CoMFAに実装されているPLSは、予測子の共線性を利用して情報を縮約しながら回帰式を得る方法である。3次元空間内において、近隣格子点同士の分子場ポテンシャルは類似し必然的に共線性が高くなるためPLS解析との相性がよく、各分子場ポテンシャルの重要度を評価することによって、3次元空間上でどの部位のどのような物理化学的変化が活性に影響を明らかにすることができる。これはCoMFAの手法を踏襲した本解析にも当てはまるため、今回作成したPLSモデルの結果に基づいて薬剤耐性獲得に影響する変異体構造の可視化を試みた。

Figure 9は、ロピナビルとプロテアーゼの複合体上にPLSモデルの標準化偏回帰係数(係数×標準偏差)の等高線マップをプロットしたものである。等高線の値には、標準化偏回帰係数の1パーセントおよび99パーセントを用いた。Figure 9では、緑色の等高線ではファンデルワールスポテンシャルが大きくなるほど、薬剤耐性が大きくなり、黄色の等高線ではファンデルワールスポテンシャルが小さくなるほど、薬剤耐性が大きくなる傾向を示している。黄色の等高線は、薬剤耐性に寄与するとされる残基やロピナビルとプロテアーゼの結合部分周辺に分布している。

本手法では、プロテアーゼと格子点上に配置したプローブ原子との相互作用のエネルギー(ファンデルワールスポテンシャルおよびクーロンポテンシャル)を予測子として利用している。ある格子上のファンデルワールスポテンシャルが大きくなることは、その格子点へのプロテアーゼの接近、空間的狭まりを反映し、逆に、ファンデルワールスポテンシャルが小さくなることは、プロテアーゼの離隔、空間的広がりを反映すると考えられる。したがって、Figure 9の緑色の等高線部分は、その部分の空間的な狭まりが薬剤耐性を大きくし、黄色の等高線部分は、空間的な広がりが薬剤耐性を大きくすると解釈できる。

クーロンポテンシャルについても同様の処理を行ったが、すべての標準化偏回帰係数が閾値(1パーセント、99パーセント)外に存在せず、薬剤耐性へのクーロンポテンシャルの寄与率は小さいと判断された。



**Figure 9.** Contour map of steric effects in drug resistance acquisition. Contours were generated based on PLS standardized partial regression coefficients. Yellow and green contours indicate 1st and 99th percentiles of standardized partial regression coefficients, respectively. Steric interaction of the protease with yellow regions negatively affects drug resistance acquisition, whereas green regions show a positive effect. Dimerized wild-type HIV-1 protease (grey ribbon) and lopinavir (wireframe) are shown in the same figure. Pink-colored sites of the protease indicate amino acids involved in drug resistance.

## II - 5 考察

本研究では、Stanford HIV Drug Resistance database からプロテアーゼ阻害剤に関するデータセットを構築した。このデータセットには、配列確定データと配列未確定データが含まれる。配列が未確定となる理由には、*in vitro* の HIV のサンプルに複数の HIV が混在していた、あるいはサンガー法の検出限界以下であったことが考えられる [92]。この配列未確定データに対して全てのパターンを考えることは、信頼性の低い誤りのあるデータを増加させてしまう。実際、Table 5 の通り、単純に組み合わせ配列を数え上げると、配列確定データの 100 倍以上になる。これらのデータをそのまま学習データに含めた場合には、構築するモデルが配列確定データを無視して、不確かなデータばかりに適合してしまう危険性がある。そこで、本研究では、配列未確定データから生成された組み合わせ配列の存在確率を推定し、その存在確率でデータに重みづけを行う方法を考案した。その結果、配列未確定データの実質的なサンプルサイズは配列確定データの 2 倍程度であり、配列確定データに加えて、配列未確定データが持つ情報量を有効に利用できるようになった。

予測モデルの構築については、比較的計算コストが小さく広く使われている PLS、LGBM、RF、SVR を利用した。Table 12 の通り、機械学習法間で予測精度に大きな差はなかったが、すべてのプロテアーゼ阻害剤で LGBM が比較的高い予測精度を示した。LGBM は、近年開発された新技法であり、内部計算において複雑な木構造を取ることにより、高い予測精度を出せるとされている [64]。一方で、過学習のリスクもあることから、一概に LGBM が優れる訳ではない。あくまで本データセットに対して、LGBM の性能が良好に働いたと考えるべきであろう。

既存の予測モデルとは、解析データセットやデータ処理方法が異なるため、厳密な比較は困難であるが、本研究で提案するモデルは、これまでに遜色のないレベルでの予測精度を与えている。例えば、本研究とは異なる小さなデータセットから構築された Gene2pheno と呼ばれる予測モデルでは、プロテアーゼ阻害剤の決定係数の平均値は約 0.698 であり [55]、個々のプロテアーゼ阻害剤を見た場合でも本モデルの方がすべて高い予測精度を与えている (Table 12)。また、本研究と同じデータセットを利用した ChenHsiang らの予測モデルの決定係数の平均値は約 0.88 であり、本研究よりも高い予測精度を示している [58]。しかしながら、ChenHsiang らの研究は、配列未確定データから生成したすべてのパターンに対して、重みもつけずに配列確定データとマージして予測モデルの構築を行っており、数的重みにより未確定データへの過剰適合してしまっていると考えられる。このように正解ラベルの扱いの違いにより、必ずしも本研究での予測性能が劣るとは言えない。一方、HIV が薬剤耐性を持つか否かという分類

をする予測モデルの開発も進められている。そこで、モデルの分類性能についても評価を行った (Table 13)。近年様々な分野で成果を上げている Deep Learning を応用した過去の分類モデル [54] と比較することとし、そのモデルで評価に使われている複数のメトリクスを利用した。一般に、分類性能の評価には、分類されるデータのバランス、例えば、A と B に分類する場合、A:B = 5 : 5 のようにデータの比率のバランスがよい場合に適当とされるメトリクス (accuracy、precision、true positive ratio) と、A:B = 9 : 1 のように比率のバランスの悪い場合に適当とされるメトリクス (AUC、F1) がある。本研究での分類性能は、チプラナビルを除けば、様々なメトリクスで約 0.9 程度であり、ほぼ同等の予測精度であった。

チプラナビルへの薬剤耐性の予測精度は、他の薬剤と比較すると低かった。機械学習の性能は、トレーニングデータの質に依存するため、たとえ同じ解析法を用いても同等の予測精度が得られるとは保証できない。しかし、それだけではなく、予測精度が低い理由には、チプラナビルが他のプロテアーゼ阻害剤と違う作用機序を持つことと関係しているかもしれない。HIV プロテアーゼは、二量体を形成し、二量体の中心部分が活性部位となる [93]。一般に、プロテアーゼ阻害剤は、二量体の活性部位に結合し阻害作用を示すが、チプラナビルは、この阻害効果に加えて、二量体化自体の阻害作用も持つことが報告されている [94]。今回構築したモデルは二量体構造を対象としているので、この副次的な作用までを扱うことができない点が問題点として挙げられる。ダルナビルも、チプラナビルと同様の作用機序を持つことが知られている [94, 95] が、本研究での予測精度はおおむね良好であった。両薬剤での矛盾に関しては二つの作用機序の寄与率の違いで説明されるかもしれないが、これを明らかにするにはプロテアーゼの二量体形成に関わる影響を分離評価した定量的データが必要であり、現時点では仮説の域を超えることができない。

CoMFA 等高線マップより、薬剤周辺のスペースが広がるほど、薬剤耐性が大きくなることが示唆された。Wang ら [96] は、プロテアーゼの変異体において、薬剤とプロテアーゼとの結合部分の広がり、薬剤とプロテアーゼとの結合を弱め、結果として薬剤耐性を増大させることを分子動力学シミュレーションを用いて報告しており、本研究で得られた傾向は妥当なものと考えられる。また、緑色の等高線は、黄色の等高線の外側に対応するかのよう分布している (Figure 9)。緑色の周辺の空間的な狭まりが薬剤耐性を増大させるのは、先に述べた薬剤とプロテアーゼの結合スペース部分の広がり、と外側の空間の狭まりとが対応関係にある可能性を示唆している。

## 結論

以上、著者は二章にわたり、HIVの抗ウイルス治療を対象として、初回療法での第一選択レジメンの有効性の比較・評価およびHIV薬剤耐性の予測という階層の異なる2つの課題に対して、それぞれの局面に応じた新しいモデル解析法を提示し、以下の結論を得た。

### 第一章 モデルベースメタ解析法による初回療法での最も有効なレジメンの選択

モデルベースメタ解析法を導入して、エファビレンツベースのレジメンおよびドルテグラビルベースのレジメンの初回療法での有効性を定量的に評価した。本手法により、臨床試験設定の異なるデータセットを統合して、大規模なデータセットに基づく解析結果を得ることができた。治療効果を速度論パラメータからなる数理モデルで記述した本手法の適用により、エファビレンツベースのレジメンと比較して、ドルテグラビルベースのレジメンでは、治療効果の約4倍発現が早く、治療効果が約2倍持続しやすいことが明らかとなった。さらに、こうした速度過程ごとでの定量的な比較評価に加え、各速度論パラメータに影響する共変量の解析も行った結果、レジメンに依らず、治療開始前のウイルス量が多いあるいはCD4のカウントが少ない感染者つまり重症者ほど予後が悪い傾向があるが示唆された。

### 第二章 HIVの薬剤耐性に関する予測モデル構築法の提案と検証

HIVの薬剤耐性の予測モデルの構築法を提案するために、プロテアーゼ阻害剤を対象として予測モデルの構築、検証を行った。集団遺伝学をベースとした考えに則って、データの拡張を行い、大規模なデータセットを構築することに成功した。さらに、プロテアーゼの変異体の立体構造変化に注目し、立体構造を反映した分子場マッピングを応用して特徴量を生成し、機械学習モデルに組み込み予測モデルの構築を行った。その結果、得られた回帰の決定係数は0.51~0.86程度、分類問題に投射した際の正解率は90%程度となり、外部検証データで十分な予測精度を持つモデルを構築することができた。さらに、予測子となる特徴量の重要度を分子場上に空間マッピングしたところ、薬剤耐性に寄与する特徴量情報が、薬剤耐性に関与することが報告されているアミノ酸残基の周辺及び薬剤とプロテアーゼの結合サイト周辺に認められ、本解析方法がウイルス薬剤耐性の予測モデル開発に有効であることが確認できた。

以上、著者は、抗HIV治療に対して、初回療法での第一選択レジメンの有効性の比較・評価およびHIV薬剤耐性の予測を行った。得られた知見は、HIV治療の最適化に有

益な情報を提示するとともに、今後のレジメン評価・予測モデル開発の指針を示すものとする。

## 謝辞

終わりに臨み、本研究の実施にあたり、終始御懇篤なる御指導、御鞭撻を賜りました京都大学大学院薬学研究科 山下富義教授に衷心より深甚なる謝意を表します。

また、終始御懇切なる御指導を賜りました京都大学大学院薬学研究科 樋口ゆり子准教授に謹んで深く感謝の意を表します。

本研究に際して、無償での計算環境の提供を頂きました情報・システム研究機構国立遺伝学研究所様へ深く感謝の意を表します。

さらに、様々な有益な御指摘、御助言を賜りました京都大学大学院薬学研究科 津田真弘講師、宗可奈子助教、並びに、様々なご助言を頂きました京都大学大学院薬学研究科 薬品動態制御学分野・実践臨床薬学分野教室員一同に深謝いたします。

加えて、著者は、日本学術振興会特別研究員の制度により支援を受けましたので、ここに感謝いたします。

最後に、研究に専念できる環境を与えてくださった父 広和、母 由香に深く感謝いたします。



# 実験の部

## 第一章 実験の部

### [1] データの収集

PubMed, Cochrane Central Register of Controlled Trials (CENTRAL), and clinicaltrials.gov の三つの独立したデータベースから、ドルテグラビルおよびエファビレンツに関する臨床試験情報を収集する。PubMedは、最も複雑な検索式に基づく文献検索が可能であるが、他のデータベースに関しては、不可能であるためそれぞれ適当な検索式を採用する必要がある。そのため、Table 14のように検索式をそれぞれのデータベースに合わせて決定した。各データベースから収集された論文の重複削除には、Refworks (version 2.0)を利用した。そのうち、タイトル、要旨から明らかに適格条件に合わない論文を取り除き、タイトル・要旨からは判断ができなかった論文は、全文を読了し決定を下すことで、最終的に解析可能なデータセットを取得した。

### [2] モデルの構築

本論文中式(1) (4)のパラメータを収集された臨床試験データに対する当てはめ計算によって求めた。この際、解析には、NONMEM7.3のFOCE-Iを用いて行った。その他、統計量においてもすべて、NONMEM7.3から出力される値を利用した。

**Table 14.** 各データベースで利用した検索式

|                    | Dolutegravir  | Efavirenz   |
|--------------------|---|---|
| PubMed             | ((randomized controlled trial[pt] OR randomized controlled trials[mh] OR random allocation[mh] OR controlled clinical trial[pt] OR randomized[tw] OR randomised[tw] OR randomly[tw] OR random*[tw] OR trial[tiab] OR groups[tiab]) AND (dolutegravir[tw] OR DTG[tw] OR S/GSK1349572[tw] OR GSK-1349572[tw] OR Tivicay[tw] OR dolutegravir[Supplementary Concept]) ) | ((randomized controlled trial[pt] OR randomized controlled trials[mh] OR random allocation[mh] OR controlled clinical trial[pt] OR randomized[tw] OR randomised[tw] OR randomly[tw] OR random*[tw] OR trial[tiab] OR groups[tiab]) AND (Efavirenz[tw] OR EFV[tw] OR EFZ[tw] OR efavirenz, R-isomer[tw] OR Sustiva[tw] OR L 743726[tw] OR L-743,726[tw] OR L-743726[tw] OR L 743,726[tw] OR Stocrin[tw] OR Merck Sharp and Dohme brand of efavirenz[tw] OR DMP 266[tw] OR DMP-266[tw] OR dmp266[tw] OR efavir[tw] OR lginase[tw] OR l743726[tw] OR efavirenz, S-isomer[tw] OR virorrever[tw])) |
| CENTRAL            | (randomized controlled trial OR randomized controlled trials OR random allocation OR controlled clinical trial OR randomized OR randomised OR randomly OR random OR trial OR groups) AND (dolutegravir OR DTG OR S/GSK1349572 OR GSK-1349572 OR Tivicay)  | ((randomized controlled trial OR randomized controlled trials OR random allocation OR controlled clinical trial OR randomized OR randomised OR randomly OR random OR trial OR groups) AND (Efavirenz OR EFV OR EFZ OR efavirenz, R-isomer OR Sustiva OR L 743726 OR L-743,726 OR L-743726 OR L 743,726 OR Stocrin OR Merck Sharp and Dohme brand of efavirenz OR DMP 266 OR DMP-266 OR dmp266 OR efavir OR lginase OR l743726 OR efavirenz, S-isomer OR virorrever))  |
| clinicaltrials.gov | dolutegravir OR DTG OR S/GSK1349572 OR GSK-1349572 OR Tivicay OR dolutegravir   Studies With Results   Interventional Studies   | Efavirenz OR EFV OR EFZ OR efavirenz, R-isomer OR Sustiva OR L 743726 OR L-743,726 OR L-743726 OR L 743,726 OR Stocrin OR Merck Sharp and Dohme brand of efavirenz OR DMP 266 OR DMP-266 OR dmp266 OR efavir OR lginase OR l743726 OR efavirenz,S-isomer   Studies With Results   Interventional Studies  |

## 第二章 実験の部

### [1] データの収集

Stanford HIV Drug resistance database から、プロテアーゼの一次配列と薬剤耐性インデックスがペアになったデータを収集した。薬剤耐性インデックスとは、 $FC = \frac{\text{mutant の IC50}}{\text{wild type の IC50}}$  でデータベース上、定義されている。また、従来の報告と同様に、プロテアーゼにおいては、 $FC \geq 3.5$  を薬剤耐性あり、 $FC < 3.5$  を薬剤耐性なしとして分類した。

### [2] データの前処理

ホモロジーモデリングには、Modeller (sci 9.20)[97] を利用した。ホモロジーモデリングに用いる鋳型のプロテアーゼの構造は、Protein Data Bank から得た。その際、薬剤が結合した共結晶であり、結晶化の実験条件が同じという制約をかけて選択した。予測された立体構造上の各原子のチャージは、PDB2PQR (version 2.1.1)[98] を用いて、pH=7.0 という条件下で計算した。

### [3] モデルの構築

Recursive feature elimination (RFE) には、scikit-learn (version 0.21.3) を利用した。ハイパーパラメータの最適化には、Optuna (version 2.0.0)[88] を利用した。機械学習の各種モデルについては、LightGBM には、LightGBM(version 2.3.0) の scikit-learn API を利用し、Support vector regression および Random forest regression には、scikit-learn (version 0.21.3) を利用し、Partial least squares は、申請者が Python 3.7 上で実装した。

### [4] 薬剤耐性ウイルスの構造要因解析

等高線マップの作成には、従来の CoMFA のアプローチと同様に、各特徴量の空間座標に対応する PLS モデルの標準化偏回帰係数(係数×標準偏差)を利用した。まず、特徴抽出により残された空間座標同士をドロネー四面体で分割し、メッシュを作成した。次に、メッシュのすべての辺上で、標準化偏回帰係数の1パーセントおよび99パーセントの値を示す内分点がないかどうかを探索した。最後に、内分点同士を結んで、等高線マップを作成した。この際、内分点同士の距離が5 Å 以上の場合は、別のクラスターとみなして、空間上に等高線をマッピングした。図の作成には、Chimera (version 1.15)[99] を利用した。

## 引用文献

1. UNAIDS data 2020. [https://www.unaids.org/sites/default/files/media\\_asset/2020\\_aids-data-book\\_en.pdf](https://www.unaids.org/sites/default/files/media_asset/2020_aids-data-book_en.pdf). Accessed 27 Jan 2021.
2. 90-90-90: good progress, but the world is off-track for hitting the 2020 targets | UNAIDS. [https://www.unaids.org/en/resources/presscentre/featurestories/2020/september/20200921\\_90-90-90](https://www.unaids.org/en/resources/presscentre/featurestories/2020/september/20200921_90-90-90). Accessed 27 Jan 2021.
3. Nachega JB, Marconi VC, Zyl GU van, Gardner EM, Preiser W, Hong SY, et al. HIV treatment adherence, drug resistance, virologic failure: Evolving concepts. *Infect Disord - Drug Targets*. 2011;11:167-74. doi:10.2174/187152611795589663.
4. Lee FJ, Amin J, Carr A. Efficacy of initial antiretroviral therapy for HIV-1 infection in adults: A systematic review and meta-analysis of 114 studies with up to 144 weeks' follow-up. *Plos One*. 2014;9.
5. Maggiolo F. Efavirenz: A decade of clinical experience in the treatment of HIV. *J Antimicrob Chemoth*. 2009;64:910-928.
6. Consolidated guidelines on the use of antiretroviral drugs for treating and preventing HIV infection. [https://apps.who.int/iris/bitstream/handle/10665/198064/9789241509893\\_eng.pdf?sequence=1](https://apps.who.int/iris/bitstream/handle/10665/198064/9789241509893_eng.pdf?sequence=1). Accessed 27 Jan 2021.
7. Clinton health access initiative. HIV market report. [https://clintonhealthaccess.org/content/uploads/2018/09/2018-HIV-Market-Report\\_FINAL.pdf](https://clintonhealthaccess.org/content/uploads/2018/09/2018-HIV-Market-Report_FINAL.pdf). Accessed 27 Jan 2021.
8. Updated recommendations on first-line and second-line antiretroviral regimens and post-exposure prophylaxis and recommendations on early infant diagnosis of HIV: interim guidance. Updated recommendations on first-line and second-line antiretroviral regimens and post-exposure prophylaxis and recommendations on early infant diagnosis of HIV: interim guidance. 2018.
9. Rutherford GW, Horvath H. Dolutegravir plus two nucleoside reverse transcriptase inhibitors versus efavirenz plus two nucleoside reverse transcriptase inhibitors as initial antiretroviral therapy for people with HIV: A systematic review. *Plos One*. 2016;11.
10. Gross JL, Rogers J, Polhamus D, Gillespie W, Friedrich C, Gong Y, et al. A novel model-based meta-analysis to indirectly estimate the comparative efficacy of two medications: An example using DPP-4 inhibitors, sitagliptin and linagliptin, in treatment of type 2 diabetes mellitus. *Bmj Open*. 2013;3.
11. Mould DR. Model-based meta-analysis: An important tool for making quantitative decisions during drug development. *Clin Pharmacol Ther*. 2012;92:283-6.

12. Kryst J, Kawalec P, Pilc A. Efavirenz-based regimens in antiretroviral-naive HIV-infected patients: A systematic review and meta-analysis of randomized controlled trials. *Plos One*. 2015;10.
13. Jpt H, Green S. *Cochrane Handbook for Systematic Reviews of Interventions* 5.0.1. 2008.
14. Moher D, Liberati A, Tetzla J, Altman DG, Altman D, Antes G, et al. Preferred reporting items for systematic reviews and meta-analyses: The PRISMA statement. *Plos Med*. 2009;6.
15. Walmsley S, Baumgarten A, Berenguer J, Felizarta F, Florence E, Khuong-Josses M-A, et al. Dolutegravir plus abacavir/lamivudine for the treatment of HIV-1 infection in antiretroviral therapy-naive patients: Week 96 and week 144 results from the SINGLE randomized clinical trial. *J Aids J Acquir Immune De c Syndromes*. 2015;70:515-9.
16. Sax PE, Tierney C, Collier AC, Fischl MA, Mollan K, Peeples L, et al. Abacavir-lamivudine versus tenofovir-emtricitabine for initial HIV-1 therapy. *New Engl J Medicine*. 2009;361:2230-40.
17. Honda M, Ishisaka M, Ishizuka N, Kimura S, Oka S. Open-label randomized multicenter selection study of once daily antiretroviral treatment regimen comparing ritonavir-boosted atazanavir to efavirenz with fixed-dose abacavir and lamivudine. *Internal Med*. 2011;50:699-705.
18. Echeverria P, Negredo E, Carosi G, Galvez J, Gomez JL, Ocampo A, et al. Similar antiviral efficacy and tolerability between efavirenz and lopinavir/ritonavir, administered with abacavir/lamivudine (Kivexa), in antiretroviral-naive patients: A 48-week, multicentre, randomized study (Lake Study). *Antivir Res*. 2010;85:403-8.
19. Post FA, Moyle GJ, Stellbrink HJ, Domingo P, Podzamczar D, Fisher M, et al. Randomized comparison of renal effects, efficacy, and safety with once-daily abacavir/lamivudine versus tenofovir/emtricitabine, administered with efavirenz, in antiretroviral-naive, HIV-1-infected adults: 48-week results from the ASSERT study. *J Aids J Acquir Immune De c Syndromes*. 2010;55:49-57.
20. Kumar P, DeJesus E, Huhn G, Sloan L, Small CB, Edelstein H, et al. Evaluation of cardiovascular biomarkers in a randomized trial of fosamprenavir/ritonavir vs. efavirenz with abacavir/lamivudine in underrepresented, antiretroviral-naive, HIV-infected patients (SUPPORT): 96-week results. *Bmc Infect Dis*. 2013;13.
21. Bartlett JA, Johnson J, Herrera G, Sosa N, Rodriguez A, Liao Q, et al. Long-term results of initial therapy with abacavir and lamivudine combined with efavirenz, amprenavir/ritonavir, or stavudine. *J Aids J Acquir Immune De c Syndromes*. 2006;43:284-92.
22. Moyle GJ, DeJesus E, Cahn P, Castillo SA, Zhao H, Gordon DN, et al. Abacavir once or twice daily combined with once-daily lamivudine and efavirenz for the treatment of antiretroviral-naive HIV-infected adults: Results of the ziagen once daily in antiretroviral combination study. *J Aids J Acquir Immune De c Syndromes*. 2005;38:417-25.
23. DeJesus E, Herrera G, Teo lo E, Gerstoft J, Buendia CB, Brand JD, et al. Abacavir versus zidovu-

dine combined with lamivudine and efavirenz, for the treatment of antiretroviral-naive HIV-infected adults. *Clin Infect Dis*. 2004;39:1038-46.

24. Podzamczer D, Ferrer E, Sanchez P, Gatell JM, Crespo M, Fisac C, et al. Less lipoatrophy and better lipid profile with abacavir as compared to stavudine: 96-Week results of a randomized study. *J Acquir Immune Defic Syndromes*. 2007;44:139-47.

25. Kravchenko A, Orlova-Morozova E, Nagimova F, Kozirev O, Shimonova T, Bichko V. Safety and antiviral effect of Elpida (VM-1500), a novel NNRTI (+Truvada) in treatment-naive HIV-1-infected patients at 24- to 48-week therapy. *J Int AIDS Soc*. 2016;19:36.

26. Arribas JR, Pozniak AL, Gallant JE, DeJesus E, Gazzard B, Campo RE, et al. Tenofovir disoproxil fumarate, emtricitabine, and efavirenz compared with zidovudine/lamivudine and efavirenz in treatment-naive patients: 144-Week analysis. *J Acquir Immune Defic Syndromes*. 2008;47:74-8.

27. Rockstroh JK, DeJesus E, Lennox JL, Yazdanpanah Y, Saag MS, Wan H, et al. Durable efficacy and safety of raltegravir versus efavirenz when combined with tenofovir/emtricitabine in treatment-naive HIV-1-infected patients: Final 5-year results from STARTMRK. *J Acquir Immune Defic Syndromes*. 2013;63:77-85.

28. Molina J-M, Cahn P, Grinsztejn B, Lazzarin A, Mills A, Saag M, et al. Rilpivirine versus efavirenz with tenofovir and emtricitabine in treatment-naive adults infected with HIV-1 (ECHO): A phase 3 randomised double-blind active-controlled trial. *Lancet*. 2011;378:238-46.

29. Landman R, Koulla-Shiro S, Sow PS, Ngolle M, Diallo M-B, Gueye NFN, et al. Evaluation of four tenofovir-containing regimens as first-line treatments in Cameroon and Senegal: The ANRS 12115 DAYANA trial. *Antivir Ther*. 2014;19:51-9.

30. Vernazza P, Wang C, Pozniak A, Weil E, Pulik P, Cooper DA, et al. Efficacy and safety of ledipasvirine (UK-453,061) versus efavirenz in antiretroviral treatment-naive HIV-1-infected patients: Week 48 primary analysis results from an ongoing, multicenter, randomized, double-blind, phase IIb trial. *J Acquir Immune Defic Syndromes*. 2013;62:171-9.

31. Cohen C, Elion R, Ruane P, Shamblaw D, DeJesus E, Rashbaum B, et al. Randomized, phase 2 evaluation of two single-tablet regimens elvitegravir/cobicistat/emtricitabine/tenofovir disoproxil fumarate versus efavirenz/emtricitabine/tenofovir disoproxil fumarate for the initial treatment of HIV infection. *Aids*. 2011;25:F7-12.

32. Amin J, Becker S, Belloso W, Boito M, Cooper D, Crabtree-Ramirez B, et al. Efficacy and safety of efavirenz 400 mg daily versus 600 mg daily: 96-week data from the randomised, double-blind, placebo-controlled, non-inferiority ENCORE1 study. *Lancet Infect Dis*. 2015;15:793-802.

33. Amin J, Becker S, Belloso W, Boito M, Cooper D, Crabtree-Ramirez B, et al. Efficacy of 400 mg efavirenz versus standard 600 mg dose in HIV-infected, antiretroviral-naive adults (ENCORE1):

- A randomised, double-blind, placebo-controlled, non-inferiority trial. *Lancet*. 2014;383:1474-82.
34. Sax PE, DeJesus E, Mills A, Zolopa A, Cohen C, Wohl D, et al. Co-formulated elvitegravir, cobicistat, emtricitabine, and tenofovir versus co-formulated efavirenz, emtricitabine, and tenofovir for initial treatment of HIV-1 infection: A randomised, double-blind, phase 3 trial, analysis of results after 48 weeks. *Lancet*. 2012;379:2439-48.
  35. Lunzen JV, Antinori A, Cohen CJ, Arribas JR, Wohl DA, Rieger A, et al. Rilpivirine vs. efavirenz-based single-tablet regimens in treatment-naive adults: Week 96 efficacy and safety from a randomized phase 3b study. *Aids*. 2016;30:251-9.
  36. Thompson M, Saag M, Dejesus E, Gathe J, Lalezari J, Landay AL, et al. A 48-week randomized phase 2b study evaluating cenicriviroc versus efavirenz in treatment-naive HIV-infected adults with C-C chemokine receptor type 5-tropic virus. *Aids*. 2016;30:869-78.
  37. Miro JM, Manzardo C, Ferrer E, Lonca M, Guardo AC, Podzamczer D, et al. Immune reconstitution in severely immunosuppressed antiretroviral-naive HIV-1-infected patients starting efavirenz, lopinavir-ritonavir, or atazanavir-ritonavir plus tenofovir/emtricitabine: Final 48-week results (The Advanz-3 Trial). *J Acquir Immune Defic Syndromes*. 2015;69:206-15.
  38. Puls RL, Srasuebku P, Petoumenos K, Boesecke C, Duncombe C, Belloso WH, et al. Efavirenz versus boosted atazanavir or zidovudine and abacavir in antiretroviral treatment - Naive, HIV-infected subjects: Week 48 data from the Altair study. *Clin Infect Dis*. 2010;51:855-64.
  39. Markowitz M, Nguyen B-Y, Gotuzzo E, Mendo F, Ratanasuwan W, Kovacs C, et al. Sustained antiretroviral effect of raltegravir after 96 weeks of combination therapy in treatment-naive patients with HIV-1 infection. *J Acquir Immune Defic Syndromes*. 2009;52:350-6.
  40. Elvucitabine/Efavirenz/Tenofovir vs. Lamivudine/Efavirenz/Tenofovir in HIV-1 infected, treatment naive subjects. 2016.
  41. Gallant JE, Staszewski S, Pozniak AL, DeJesus E, Suleiman JMAH, Miller MD, et al. Efficacy and safety of tenofovir DF vs stavudine in combination therapy in antiretroviral-naive patients: A 3-year randomized trial. *Jama*. 2004;292:191-201.
  42. Raff F, Rachlis A, Stellbrink H-J, Hardy WD, Torti C, Orkin C, et al. Once-daily dolutegravir versus raltegravir in antiretroviral-naive adults with HIV-1 infection: 48 week results from the randomised, double-blind, non-inferiority SPRING-2 study. *Lancet*. 2013;381:735-43.
  43. Raff F, Jaeger H, Quiros-Roldan E, Albrecht H, Belonosova E, Gatell JM, et al. Once-daily dolutegravir versus twice-daily raltegravir in antiretroviral-naive adults with HIV-1 infection (SPRING-2 study): 96 week results from a randomised, double-blind, non-inferiority trial. *Lancet Infect Dis*. 2013;13:927-35.
  44. Gallant J, Lazzarin A, Mills A, Orkin C, Podzamczer D, Tebas P, et al. Bictegravir, emtric-

- itabine, and tenofovir alafenamide versus dolutegravir, abacavir, and lamivudine for initial treatment of HIV-1 infection (GS-US-380-1489): a double-blind, multicentre, phase 3, randomised controlled non-inferiority trial. *Lancet*. 2017;390:2063-72.
45. Orrell C, Hagins DP, Belonosova E, Porteiro N, Walmsley S, FalcoV, et al. Fixed-dose combination dolutegravir, abacavir, and lamivudine versus ritonavir-boosted atazanavir plus tenofovir disoproxil fumarate and emtricitabine in previously untreated women with HIV-1 infection (ARIA): week 48 results from a randomised, open-label, non-inferiority, phase 3b study. *Lancet Hiv*. 2017;4:e536-46.
46. Donahue DA, Sloan RD, Kuhl BD, Bar-Magen T, Schader SM, Wainberg MA. Stage-dependent inhibition of HIV-1 replication by antiretroviral drugs in cell culture. *Antimicrob Agents Ch*. 2010;54:1047-54.
47. Canducci F, Ceresola ER, Boeri E, Spagnuolo V, Cossarini F, Castagna A, et al. Cross-resistance profile of the novel integrase inhibitor dolutegravir (S/GSK1349572) using clonal viral variants selected in patients failing raltegravir. *J Infect Dis*. 2011;204:1811-5.
48. Tang MW, Shafer RW. HIV-1 antiretroviral resistance: Scientific principles and clinical applications. *Drugs*. 2012;72:e1-25.
49. Sluis-Cremer N, Tachedjian G. Mechanisms of inhibition of HIV replication by non-nucleoside reverse transcriptase inhibitors. *Virus Res*. 2008;134:147-56.
50. Patel DA, Snedecor SJ, Tang WY, Sudharshan L, Lim JW, Cui R, et al. 48-Week efficacy and safety of dolutegravir relative to commonly used third agents in treatment-naive HIV-1-Infected patients: A systematic review and network meta-analysis. *Plos One*. 2014;9.
51. Skowron G, Street JC, Obee EM. Baseline CD4+ cell count, not viral load, correlates with virologic suppression induced by potent antiretroviral therapy. *Aids J Acquir Immune Defic Syndromes*. 2001;28:313-9.
52. Kuritzkes DR. Drug resistance in HIV-1. *Curr Opin Virol*. 2011;1:582-9.
53. Gunthard HF, Calvez V, Paredes R, Pillay D, Shafer RW, Wensing AM, et al. Human Immunodeficiency Virus Drug Resistance: 2018 Recommendations of the International Antiviral Society-USA Panel. *Clin Infect Dis*. 2019;68:177-87.
54. Steiner MC, Gibson KM, Crandall KA. Drug Resistance Prediction Using Deep Learning Techniques on HIV-1 Sequence Data. *Viruses*. 2020;12:560.
55. Beerenwinkel N, Da-umer M, Oette M, Korn K, Hofmann D, Kaiser R, et al. Geno2pheno: estimating phenotypic drug resistance from HIV-1 genotypes. *Nucleic Acids Res*. 2003;31:3850-5.
56. Tarasova O, Biziukova N, Filimonov D, Poroikov V. A Computational Approach for the Prediction of HIV Resistance Based on Amino Acid and Nucleotide Descriptors. *Molecules*. 2018;23:2751.
57. Yu X, Weber IT, Harrison RW. Prediction of HIV drug resistance from genotype with encoded



- three-dimensional protein structure. *Bmc Genomics*. 2014;15 Suppl 5:S1.
58. Shen C, Yu X, Harrison RW, Weber IT. Automated prediction of HIV drug resistance from genotype data. *Bmc Bioinformatics*. 2016;17 Suppl 8:278.
59. Rhee S-Y, Gonzales MJ, Kantor R, Betts BJ, Ravela J, Shafer RW. Human immunodeficiency virus reverse transcriptase and protease sequence database. *Nucleic Acids Res*. 2003;31:298-303.
60. Dorn M, Silva MB e, Buriol LS, Lamb LC. Three-dimensional protein structure prediction: Methods and computational strategies. *Comput Biol Chem*. 2014;53:251-76.
61. Cramer RD, Patterson DE, Bunce JD. Comparative molecular field analysis (CoMFA). 1. Effect of shape on binding of steroids to carrier proteins. *J Am Chem Soc*. 1988;110:5959-67.
62. Lindgren F, Geladi P, Wold S. The kernel algorithm for PLS. *J Chemometr*. 1993;7:45-59.
63. Breiman L. Random Forests. *Mach Learn*. 2001;45:5-32.
64. Ke G, Meng Q, Finley T, Wang T, Chen W, Ma W, et al. LightGBM: A Highly Efficient Gradient Boosting Decision Tree. *Advances in Neural Information Processing Systems*. 2017;:3147-55.
65. Smola AJ, Schölkopf B. A tutorial on support vector regression. *Stat Comput*. 2004;14:199-222.
66. Genotype-phenotype datasets. <https://hivdb.stanford.edu/pages/genopheno.dataset.html>. Accessed 27 Jan 2021.
67. Beerenwinkel N, Schmidt B, Walter H, Kaiser R, Lengauer T, Hoffmann D, et al. Diversity and complexity of HIV-1 drug resistance: A bioinformatics approach to predicting phenotype from genotype. *Proc National Acad Sci*. 2002;99:8271-6.
68. Rogers AR, Hu C. Linkage Disequilibrium Between Loci With Unknown Phase. *Genetics*. 2009;182:839-44.
69. Strehl A, Ghosh J. Cluster Ensembles - A Knowledge Reuse Framework for Combining Multiple Partitions. *Journal of Machine Learning Research*. 2002.
70. Pethel SD, Hahs DW. Exact Test of Independence Using Mutual Information. *Entropy*. 2014;16:2839-49.
71. Mart-Renom MA, Stuart AC, Fiser A, Sanchez R, Melo F, Sali A. Comparative protein structure modeling of genes and genomes. *Annu Rev Bioph Biom*. 2000;29:291-325.
72. Sun J, Yu EY, Yang Y, Confer LA, Sun SH, Wan K, et al. Stn1-Ten1 is an Rpa2-Rpa3-like complex at telomeres. *Gene Dev*. 2009;23:2900-14.
73. King NM, Prabu-Jeyabalan M, Bandaranayake RM, Nalam MNL, Nalivaika EA, Ozen A, et al. Extreme Entropy-Enthalpy Compensation in a Drug-Resistant Variant of HIV-1 Protease. *ACS Chem Biol*. 2012;7:1536-46.
74. Ishima R, Gong Q, Tie Y, Weber IT, Louis JM. Highly conserved glycine 86 and arginine 87 residues contribute differently to the structure and activity of the mature HIV - 1 protease. *Proteins*

Struct Funct Bioinform. 2010;78:1015-25.

75. Shen C, Wang Y, Kovalevsky AY, Harrison RW, Weber IT. Amprenavir complexes with HIV - 1 protease and its drug - resistant mutants altering hydrophobic clusters. *Febs J.* 2010;277:3699-714.

76. Mahalingam B, Wang Y, Boross PI, Tozser J, Louis JM, Harrison RW, et al. Crystal structures of HIV protease V82A and L90M mutants reveal changes in the indinavir - binding site. *Eur J Biochem.* 2004;271:1516-24.

77. Chen Z, Li Y, Chen E, Hall DL, Darke PL, Culberson C, et al. Crystal structure at 1.9-A resolution of human immunode ciency virus (HIV) II protease complexed with L-735,524, an orally bioavailable inhibitor of the HIV proteases. *J Biological Chem.* 1994;269:26344-8.

78. Clemente JC, Moose RE, Hemrajani R, Whitford LRS, Govindasamy L, Reutzel R, et al. Comparing the Accumulation of Active- and Nonactive-Site Mutations in the HIV-1 Protease †. *Biochemistry-us.* 2004;43:12141-51.

79. Wong-Sam A, Wang Y-F, Zhang Y, Ghosh AK, Harrison RW, Weber IT. Drug Resistance Mutation L76V Alters Nonpolar Interactions at the Flap-Core Interface of HIV-1 Protease. *Acs Omega.* 2018;3:12132-40.

80. Stoll V, Qin W, Stewart KD, Jakob C, Park C, Walter K, et al. X-ray crystallographic structure of ABT-378 (Lopinavir) bound to HIV-1 protease. *Bioorgan Med Chem.* 2002;10:2803-6.

81. Koz sek M, Bray J, Rezacova P, Saskova K, Brynda J, Pokorna J, et al. Molecular Analysis of the HIV-1 Resistance Development: Enzymatic Activities, Crystal Structures, and Thermodynamics of Nel navir-resistant HIV Protease Mutants. *J Mol Biol.* 2007;374:1005-16.

82. Serganov A, Huang L, Patel DJ. Structural insights into amino acid binding and gene control by a lysine riboswitch. *Nature.* 2008;455:1263-7.

83. Muzammil S, Armstrong AA, Kang LW, Jakalian A, Bonneau PR, Schmelmer V, et al. Unique Thermodynamic Response of Tipranavir to Human Immunode ciency Virus Type 1 Protease Drug Resistance Mutations-. *J Virol.* 2007;81:5144-54.

84. Vinter JG, Davis A, Saunders MR. Strategic approaches to drug design. I. An integrated software framework for molecular modelling. *J Comput Aid Mol Des.* 1987;1:31-51.

85. Wang X-S. Derivation of Coulomb ' s Law of Forces Between Static Electric Charges Based on Spherical Source and Sink Model of Particles. *Arxiv.* 2006.

86. Awasthi M, Singh S, Pandey VP, Dwivedi UN. CoMFA and CoMSIA-based designing of resveratrol derivatives as amyloid-beta aggregation inhibitors against Alzheimer ' s disease. *Med Chem Res.* 2018;27:1167-85.

87. Guyon I, Weston J, Barnhill S, Vapnik V. Gene Selection for Cancer Classi cation using Support Vector Machines. *Mach Learn.* 2002;46:389-422.

88. Teredesai A, Kumar V, Li Y, Rosales R, Terzi E, Karypis G, et al. Optuna: A Next-generation Hyperparameter Optimization Framework. *Applied Data Science Track Paper*. 2019;:2623-31.
89. Kohavi R. A Study of Cross Validation and Bootstrap for Accuracy Estimation and Model Selection. *Proceedings of the 14th international joint conference on Artificial intelligence*. 1995.
90. Fawcett T. An introduction to ROC analysis. *Pattern Recogn Lett*. 2006;27:861-74.
91. Chicco D, Jurman G. The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. *Bmc Genomics*. 2020;21:6.
92. Franca LTC, Carrilho E, Kist TBL. A review of DNA sequencing techniques. *Q Rev Biophys*. 2002;35:169-200.
93. Todd MJ, Semo N, Freire E. The structural stability of the HIV-1 protease<sup>11</sup>Edited by P. E. Wright. *J Mol Biol*. 1998;283:475-88.
94. Koh Y, Aoki M, Danish ML, Aoki-Ogata H, Amano M, Das D, et al. Loss of Protease Dimerization Inhibition Activity of Darunavir Is Associated with the Acquisition of Resistance to Darunavir by HIV-1. *J Virol*. 2011;85:10079-89.
95. Hayashi H, Takamune N, Nirasawa T, Aoki M, Morishita Y, Das D, et al. Dimerization of HIV-1 protease occurs through two steps relating to the mechanism of protease dimerization inhibition by darunavir. *Proc National Acad Sci*. 2014;111:12234-9.
96. Wang R-G, Zhang H-X, Zheng Q-C. Revealing the binding and drug resistance mechanism of amprenavir, indinavir, ritonavir, and nelfinavir complexed with HIV-1 protease due to double mutations G48T/L89M by molecular dynamics simulations and free energy analyses. *Phys Chem Chem Phys*. 2020;22:4464-80.
97. Fiser A, Sali A. Modeller: Generation and Refinement of Homology-Based Protein Structure Models. *Methods Enzymol*. 2003;374:461-91.
98. Dolinsky TJ, Nielsen JE, McCammon JA, Baker NA. PDB2PQR: an automated pipeline for the setup of Poisson-Boltzmann electrostatics calculations. *Nucleic Acids Res*. 2004;32 suppl\_2:W665-7.
99. Pettersen EF, Goddard TD, Huang CC, Couch GS, Greenblatt DM, Meng EC, et al. UCSF Chimera-A visualization system for exploratory research and analysis. *J Comput Chem*. 2004;25:1605-12.