

A Bayesian Reinforcement Learning Framework for Optimizing Sequential Combination Antiretroviral Therapy in People with HIV

Yanxun Xu

Department of Applied Math and Statistics Mathematics Institute for Data Science Division of Biostatistics and Bioinformatics The Sidney Kimmel Comprehensive Cancer Center Johns Hopkins University

> September 17, 2022 PCIC

Sequential Decision Making in HIV



Drug classes:

- Nucleotide reverse transcriptase inhibitor (NRTI)
- Non-nucleotide reverse transcriptase inhibitor (NNRTI)
- Protease inhibitor (PI)
- integrase inhibitor (INSTI)
- Entry inhibitor (EI)
- Pharmacokinetic enhancer (Booster)

A combination of ART drugs





Precision Medicine in HIV

Semi-annual visits



- Sociodemographics
- Laboratory tests
- Behavioral characteristics
- Clinical factors

Assign ART regimen



Precision Medicine in HIV

CLINICAL INFO							Enter Your Search Term	\$ Q
								In affiliation with <u>HIV.go</u>
Guidelines	Drug Database	Glossary	News	Mobile Application	Resources	Contact Us		Language (EN) -
HOME > GU	IDELINES > Initia	tion of Antire	troviral Th	erapy				
Guid	lelines	for	the	Use of A	Antir	etrov	viral Agents i	in Adults

The information in the brief version is excerpted directly from the full-text guidelines. The brief version is a compilation of the tables and boxed recommendations.

and Adolescents Living with HIV

Guidelines search Guideline Search	h Term		pen 🔻			
Version: BRIEF FULL	Initiation of Antiretroviral Th	erapy				
What's New in the Guidelines?	Updated: Dec. 18, 2019 Reviewed: Dec. 18, 2019					
Panel Roster	Panel's Recommendations for Initiating Antiretroviral Therapy in Treatment-Naive Patients					
Financial Disclosure	Panel's Recommendations					

Treatment-naive:

• Pre-treated:

change-in-sleep-habits dry-mouth upper-respiratory-illness restlessness flu-symptoms vomiting sweating sore-throat rash indigestion feeling-tired kidnev-fai sexual-problems drooling nat skin-reacti pressure -appetite² feeling-nervous nose-bl heart-failure unusual-dreams

Goal: determine the personalized ART regimen to optimize the *long-term* health

Large-Scale HIV Studies

MACS/WIHS COMBINED COHORT STUDY

• Semi-annual visits





Challenges

- Estimate the effects of ART regimens from a high-dimensional and unbalanced space.
 - High-dimensional: with more than 30 ART drugs on the market, there are a large number of possible drug combinations.
 - Unbalanced: some drug combinations are frequent whereas others are rare.



- Generate a realistic ART regimen from a large discrete space.
 - How to represent an ART regimen?

A binary vector (2^N dimension)? (Drug 1, Drug 2, ..., Drug N)



Optimize treatments from observational data (*distribution shift*).

Approach Overview



Policy Optimizer

Problem Formulation

- Baseline covariates (e.g., race): X_{i0} .
- Visit times $\boldsymbol{t}_i = (t_{i1}, \ldots, t_{i,J_i}).$
- *M* time-varying health *state* variables: $Y_i = (Y_{i1}, \ldots, Y_{i,J_i})$ with $Y_{ij} \in \mathbb{R}^M$.
- $\mathbf{Z}_i = (Z_{i1}, \ldots, Z_{i,J_i})$ with Z_{ij} denoting the cART regimen used by individual *i* during the time period $(t_{i,j-1}, t_{ij}]$.

$$\mathcal{D} = \{\mathcal{D}_i\}_{i=1}^I = \{m{X}_{i0}, m{t}_i, m{Y}_i, m{Z}_i\}_{i=1}^I$$

- State and treatment histories: $\overline{Y_{ij}} = \{Y_{ij'} : j' \leq j\}, \overline{Z_{ij}} = \{Z_{ij'} : j' \leq j\}.$
- Dynamic model: $Y_{i,j+1} = f(\overline{Y_{ij}}, \overline{Z_{i,j+1}}; \phi)$

Goal: optimize ART assignments to maximize the individual's long-term health outcomes

An Optimization Problem

- For any individual i, suppose that she already has J_i visits.
- Let $\mathbf{Y}_{i}^{\text{new}} = {\mathbf{Y}_{ij} : j > J_i}$ and $\mathbf{Z}_{i}^{\text{new}} = {Z_{ij} : j > J_i}$ denote her future longitudinal states and ART regimens.
- ART policy function: $\pi(Z_{i,j+1} \mid \overline{Y_{ij}}, \overline{Z_{ij}}; \theta)$.
- A stochastic reward function: $r_i(\mathbf{Y}_i^{\text{new}})$.

Denote the expected reward for any individual i to be:

$$R_{i}(\boldsymbol{\theta}) = \int E[\mathbf{Y}_{i}^{\text{new}}, \mathbf{Z}_{i}^{\text{new}}) \sim p(\mathbf{Y}_{i}^{\text{new}}, \mathbf{Z}_{i}^{\text{new}} | \mathcal{D}, \boldsymbol{\phi}, \boldsymbol{\theta})}[r_{i}(\mathbf{Y}_{i}^{\text{new}})] p(\boldsymbol{\phi} \mid \mathcal{D}) d\boldsymbol{\phi}.$$

Uncertainty

Goal: find the optimal policy function $\pi(\cdot, \cdot; \boldsymbol{\theta}_i^{\star})$ such that

$$\boldsymbol{\theta}_i^{\star} = \arg \max_{\boldsymbol{\theta}} R_i(\boldsymbol{\theta}).$$

An Optimization Problem

Policy gradient via stochastic gradient descent (SGD)

$$\boldsymbol{\theta}_{i,q+1} = \boldsymbol{\theta}_{i,q} + s_{i,q} \nabla_{\boldsymbol{\theta}} R_i(\boldsymbol{\theta}) \mid_{\boldsymbol{\theta} = \boldsymbol{\theta}_{i,q}}$$

$$R_i(\boldsymbol{\theta}) = \int E_{(\boldsymbol{Y}_i^{\text{new}}, \boldsymbol{Z}_i^{\text{new}}) \sim p(\boldsymbol{Y}_i^{\text{new}}, \boldsymbol{Z}_i^{\text{new}} | \mathcal{D}, \boldsymbol{\phi}, \boldsymbol{\theta})} [r_i(\boldsymbol{Y}_i^{\text{new}})] p(\boldsymbol{\phi} \mid \mathcal{D}) d\boldsymbol{\phi}.$$

$$\nabla_{\theta} R_{i}(\theta) = \int E_{(\mathbf{Y}_{i}^{\text{new}}, \mathbf{Z}_{i}^{\text{new}})} \left[\mathbf{Y}_{i}^{\text{new}}, \mathbf{Z}_{i}^{\text{new}} | \mathcal{D}, \phi, \theta} \right] \left[\mathbf{Y}_{i}^{\text{new}}, \mathbf{Y}_{i}^{\text{new}} \right] \nabla_{\theta} \log \left(\prod_{j \ge J_{i}} \pi(Z_{i,j+1} \mid \overline{\mathbf{Y}_{ij}}, \overline{Z_{ij}}; \theta) \right) \right] p(\phi \mid \mathcal{D}) d\phi$$
• Sample future states
• Define a reward function
• Parameterize the policy function

Modeling Longitudinal States



$$Y_{im}(t) = f_{im}(t) + \epsilon_{im}$$

 $(f_{i1}(t), \ldots, f_{iM}(t))$ are MGP-distributed

- Mean $(\mu_{i1}(t), \ldots, \mu_{iM}(t)).$
- Separable covariance function $cov(f_{im}(t), f_{im'}(t')) = C^M_{mm'}C^t(t, t').$
- C^M : $M \times M$ covariance matrix.
- Ornstein-Uhlenbeck (OU) kernel $C^{t}(t, t') = \rho_{t}^{|t-t'|}$.







ART Regimen Similarity



• Desired properties: (1) sharing-of-information; (2) reducing the high-dimensionality.

ART Regimen Similarity



Linear Kernel

- Computes the similarity between regimens based on the proportion of **common drugs** that two regimens share.

D4T (NRTI) + LAM (NRTI) + NFV (PI)

D4T + LAM + ATZ (PI)

D4T + LAM + EFV (NNRTI)

Representative ART Regimen

FTC (NRTI) + TDF (NRTI) + RTV (PI)

ART Regimen Similarity



Representative ART Regimen

Subset-tree (ST) Kernel

- Calculates the similarity between regimens across all levels of the tree representation.



(c) Regimen C

(d) Similarity score matrix

Modeling Longitudinal States

Posterior Inference:

- Assign priors to all unknown parameters
- Obtain posterior distributions from MCMC

Sample future states



- Define a reward function
- Parameterize the policy function



Define the reward function based on

viral load, kidney function, and depression in the next two years



Personalized weights: $\boldsymbol{w}_i = (w_{i1}, w_{i2}, w_{i3})$

An Uncertainty-Penalized Reward

Distribution shift: Model-based uncertainty-penalized policy optimization (Yu et al. 2020)

A pessimistic environment:
$$\widetilde{r}_i(Y_i^{\text{new}}) = r_i(Y_i^{\text{new}}) - \lambda u(Y_i^{\text{new}}, Z_i^{\text{new}})$$

Tuning parameter Uncertainty

$$u(\mathbf{Y}_{i}^{\text{new}}, \mathbf{Z}_{i}^{\text{new}}) = \sum_{j=J_{i}+1}^{J_{i}+4} \sum_{m=1}^{M} \sqrt{\operatorname{Var}(Y_{ijm} \mid Z_{ij}, \mathcal{D})}$$
 posterior predictive distribution

Sample future states

Define a reward function



Parameterize the policy function

Decision Process for Assigning ART



An Optimization Problem

Policy gradient via stochastic gradient descent (SGD)

$$\boldsymbol{\theta}_{i,q+1} = \boldsymbol{\theta}_{i,q} + s_{i,q} \nabla_{\boldsymbol{\theta}} R_i(\boldsymbol{\theta}) \mid_{\boldsymbol{\theta} = \boldsymbol{\theta}_{i,q}}$$



WIHS Data Analysis

- I=339 women from the Washington DC site.
- *M*=4 state variables at each visit: depression, viral load, eGFR, and BMI.
- •8% missing rate.
- Baseline covariates: age, smoking status, substance use, employment status, hypertension, and diabetes.
- N=31 ART drugs, K=6 drug classes, D=105 representative ART regimens

Precision Medicine



w = (1/3, 1/3, 1/3)

Precision Medicine



Optimal regimens for visits 32-35



Precision Medicine

Predictive depression scores under the estimated optimal regimens



23% Improvement

Precision Medicine: uncertainty-penalized policy





Precision Medicine: uncertainty-penalized policy

71 times

96 times

- When $\lambda = 0$, optimal sequence of regimens: FTC+ABC+EFV (two NR-TIs + one NNRTI)0
- When $\lambda = 0.05$, optimal sequence of regimens (3TC+ABC+EFV (two NRTIs + one NNRTI)
- When $\lambda = 0.1$, optimal sequence of regimens: 3TC+ABC+ATV+RTV (two NRTIs + one PI + one Booster)



