

rjaf: Regularized Joint Assignment Forest with Treatment Arm Clustering

May 23, 2024

Summary

Learning the assignment of treatments is an omnipresent problem in economics and public health. It arises, for example, from randomized controlled trials where a variety of behavioral nudges (treatments) are developed to enhance vaccination uptake against coronavirus disease 2019 (COVID-19) or influenza, especially among racially or ethnically underrepresented and socioeconomically disadvantaged populations (Milkman et al. 2021, 2022; Dai et al. 2021). Subject-specific covariates containing information regarding sociodemographics, clinical characteristics, and comorbid conditions, if available, can be harnessed to identify personalized treatment assignment schemes. The **rjaf** package provides a user-friendly implementation of the regularized joint assignment forest (RJAF) (Ladhania et al. 2023), a forest-based treatment assignment algorithm featuring greedy recursive partitioning (Athey, Tibshirani, and Wager 2019), treatment and covariate resampling in bootstrap tree aggregating (Breiman 1996), outcome residualization and regularization, and k-means treatment arm clustering (Hartigan and Wong 1979). Personalized treatment learning is achieved through optimizing a regularized empirical analogue of the expected outcome. The integration of R (R Core Team 2024) and C++ (Stroustrup 2013) substantially boosts the computational efficiency in tree partitioning and aggregating. This package is especially suitable in randomized controlled trial settings where a large number of treatment arms are present.

Statement of Need

There is an ever-growing literature in the intersection of machine learning and causal inference attempting to address the problem of optimal treatment assignment through heterogeneous treatment effect estimation (Athey and Imbens 2016; Wager and Athey 2018; Hitsch and Misra 2018; Athey, Tibshirani, and Wager 2019; Sverdrup et al. 2020; Athey and Wager 2021). Other methods focus on maximizing the benefit (empirical welfare) from treatment assignment (e.g., Kitagawa and Tetenov 2018), or the chance of assigning an individual to an optimal treatment arm (e.g., Murphy 2005; Zhou, Wang, and Zeng 2018). Most of these methods perform well with a limited number of treatment and control groups. As more arms are present, the estimation of arm-specific empirical welfare and the identification of individual-specific optimal arms become increasingly difficult. Commonly used implementations such as the multi-arm causal forest (Tibshirani et al. 2022) and random forest (Wright and Ziegler 2017) lead to significantly suboptimal assignment with insufficient levels of empirical welfare. By contrast, the RJAF yields elevated empirical welfare closer to the optimal level from the oracle assignment than the multi-arm causal forest and random forest. Despite the methodological advantage over existing approaches, the incorporation of machine learning and causal inference techniques such as recursive tree partitioning, bootstrap aggregating, and treatment arm clustering makes it challenging to implement the RJAF from scratch even for well-trained data scientists. The **rjaf** is an open-source software package in R and C++ that efficiently implements the RJAF, offering data scientists a user-friendly analytic toolbox for learning personalized treatment rules in real-world settings.

(simulation results available in Figure 1)

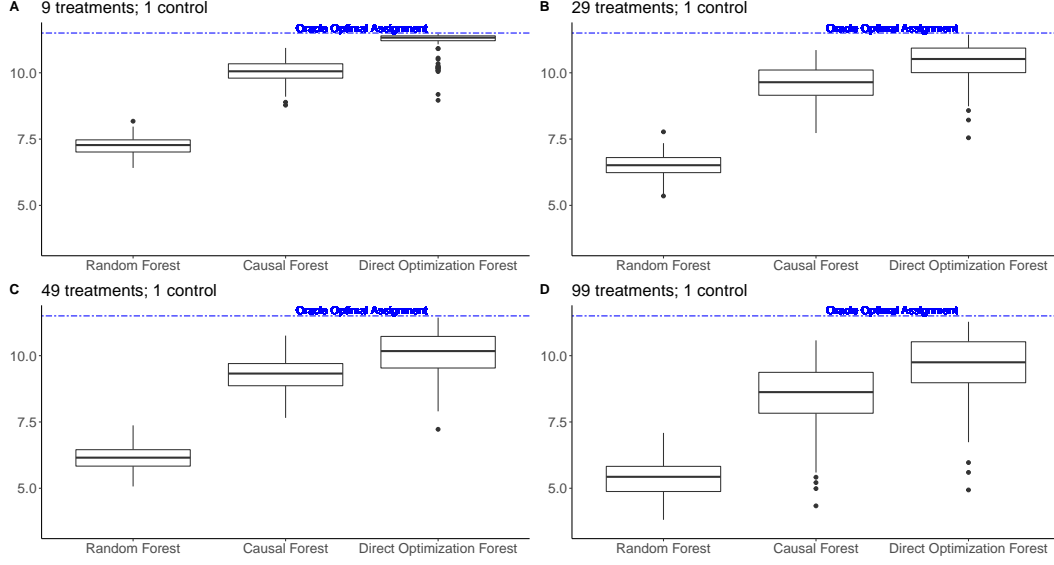


Figure 1: Direct optimization forest, random forest, and multi-arm causal forest with an increasing number of treatment arms.

Mathematics

Single dollars (\$) are required for inline mathematics e.g. $f(x) = e^{\pi/x}$

Double dollars make self-standing equations:

$$\Theta(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{else} \end{cases}$$

You can also use plain L^AT_EX for equations

$$\hat{f}(\omega) = \int_{-\infty}^{\infty} f(x)e^{i\omega x} dx \quad (1)$$

and refer to Equation 1 from text.

Acknowledgments

Wenbo Wu and Rahul Ladhania were supported by a research grant from the Robert Wood Johnson Foundation titled *Informing Strategies to Increase Use of COVID-19 and Flu Vaccines by Different Racial and Ethnic Groups to Improve Health Equity during Health Crises* (award number 78416).

References

- Athey, Susan, and Guido Imbens. 2016. “Recursive Partitioning for Heterogeneous Causal Effects.” *Proceedings of the National Academy of Sciences* 113 (27): 7353–60.
- Athey, Susan, Julie Tibshirani, and Stefan Wager. 2019. “Generalized Random Forests.” *Annals of Statistics* 47 (2): 1148–78.
- Athey, Susan, and Stefan Wager. 2021. “Policy Learning with Observational Data.” *Econometrica* 89 (1): 133–61.
- Breiman, Leo. 1996. “Bagging Predictors.” *Machine Learning* 24 (2): 123–40.

- Dai, Hengchen, Silvia Saccardo, Maria A Han, Lily Roh, Naveen Raja, Sitaram Vangala, Hardikkumar Modi, Shital Pandya, Michael Sloyan, and Daniel M Croymans. 2021. “Behavioural Nudges Increase COVID-19 Vaccinations.” *Nature* 597 (7876): 404–9.
- Hartigan, John A, and Manchek A Wong. 1979. “Algorithm AS 136: A k-Means Clustering Algorithm.” *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 28 (1): 100–108.
- Hitsch, Günter J, and Sanjog Misra. 2018. “Heterogeneous Treatment Effects and Optimal Targeting Policy Evaluation.” *Available at SSRN*.
- Kitagawa, Toru, and Aleksey Tetenov. 2018. “Who Should Be Treated? Empirical Welfare Maximization Methods for Treatment Choice.” *Econometrica* 86 (2): 591–616.
- Ladhania, Rahul, Jann Spiess, Lyle Ungar, and Wenbo Wu. 2023. “Personalized Assignment to One of Many Treatment Arms via Regularized and Clustered Joint Assignment Forests.” <https://doi.org/10.48550/arXiv.2311.00577>.
- Milkman, Katherine L, Linnea Gandhi, Mitesh S Patel, Heather N Graci, Dena M Gromet, Hung Ho, Joseph S Kay, et al. 2022. “A 680,000-Person Megastudy of Nudges to Encourage Vaccination in Pharmacies.” *Proceedings of the National Academy of Sciences* 119 (6): e2115126119.
- Milkman, Katherine L, Mitesh S Patel, Linnea Gandhi, Heather N Graci, Dena M Gromet, Hung Ho, Joseph S Kay, et al. 2021. “A Megastudy of Text-Based Nudges Encouraging Patients to Get Vaccinated at an Upcoming Doctor’s Appointment.” *Proceedings of the National Academy of Sciences* 118 (20): e2101165118.
- Murphy, Susan A. 2005. “A Generalization Error for q-Learning.” *Journal of Machine Learning Research* 6 (37): 1073–97.
- R Core Team. 2024. “R: A Language and Environment for Statistical Computing.” Vienna, Austria: R Foundation for Statistical Computing; <https://www.R-project.org>.
- Stroustrup, Bjarne. 2013. *The C++ Programming Language*. Fourth. Pearson Education.
- Sverdrup, Erik, Ayush Kanodia, Zhengyuan Zhou, Susan Athey, and Stefan Wager. 2020. “Policytree: Policy Learning via Doubly Robust Empirical Welfare Maximization over Trees.” *Journal of Open Source Software* 5 (50): 2232.
- Tibshirani, Julie, Susan Athey, Rina Friedberg, Vitor Hadad, David Hirshberg, Luke Miner, Erik Sverdrup, Stefan Wager, and Marvin Wright. 2022. “grf: Generalized Random Forests.” <https://CRAN.R-project.org/package=grf>.
- Wager, Stefan, and Susan Athey. 2018. “Estimation and Inference of Heterogeneous Treatment Effects Using Random Forests.” *Journal of the American Statistical Association* 113 (523): 1228–42.
- Wright, Marvin N., and Andreas Ziegler. 2017. “ranger: A Fast Implementation of Random Forests for High Dimensional Data in C++ and R.” *Journal of Statistical Software* 77 (1): 1–17.
- Zhou, Xuan, Yuanjia Wang, and Donglin Zeng. 2018. “Sequential Outcome-Weighted Multicategory Learning for Estimating Optimal Individualized Treatment Rules.” http://www.columbia.edu/~yw2016/SOM_ITR.pdf.