

The Feasibility and Equity of Geolocation Data for Lapse Prediction in AUD

John J. Curtin

Author note

Correspondence concerning this article should be addressed to .

Introduction

About 1 in 10 adults in the United States met diagnostic criteria for alcohol use disorder (AUD) in 2022 [1]. While some individuals will experience natural recovery (i.e., improvement without intervention) [2], for others AUD will present as a chronic, relapsing disorder marked by periods of recovery interspersed with returns back to harmful use [3], [4]. For such individuals, continued monitoring may be beneficial in assisting with the maintenance of recovery goals and in identifying precipitants to lapses, or single instances of goal-inconsistent use that may lead to relapse [5]. One sustainable and scalable way to provide this continuous monitoring to individuals who need it most is through developing algorithms to predict lapses using both personal sensing data and machine learning.

Personal sensing data are data derived via embedded sensors in technology ubiquitous in our daily lives, such as smartphones, smartwatches, or other wearables [6]. Because these devices are already so integrated within our day-to-day lives, one benefit of porting these data to clinical use is their proven ability to be collected unobtrusively and continuously. Importantly, these data do not require individuals to change their behavior or routines in any way. Moreover, when paired with machine learning models, statistical patterns connecting antecedents to lapse derived from these data (e.g., changes in mood, difficulty with close social connections, proximity to risky locations) to true lapse events can be uncovered. This is crucial for several reasons: 1) even when someone anticipates an oncoming lapse, it may be difficult to pinpoint the specific driving forces behind it; 2) these precipitating factors will have great variation both between- and within-people;

and 3) uncovering these factors may help relieve some of the cognitive burden of recovery (i.e., constant monitoring of potential environmental risk factors).

Geolocation Data for Risk Monitoring

Recovery and return to use are dynamic processes. Factors that contribute both to maintenance of recovery and return to use change from person-to-person and from moment-to-moment. A shift in social supports (e.g., a move, a break-up) may precede a lapse for one individual but not another. Time spent in locations where alcohol is available (e.g., bars, restaurants, concert venues) may precede a given lapse for another individual, but will not necessarily precede future lapses in that same individual. In order to best capture this fluidity, the ideal data type used within continuous risk monitoring systems should be able to provide a correspondingly appropriate level of granularity. One promising data source is geolocation data.

Geolocation data consist of latitude and longitude coordinates and can be sampled at regular intervals using applications on smartphones with little to no input from the user beyond initial set-up. Many smartphones and smartwatches automatically collect these data by default. This fact, paired with increasing rates of smartphone ownership, suggest that there is high potential for these data to be feasibly harnessed for use in a risk-monitoring system [7]. The importance of location, such as environmental cues or one's perceived riskiness of a setting, has been shown to play an important role in lapse [8]–[10]. This link with lapse risk has translated into the integration of coping skills that target substance-associated contexts in several treatment strategies like mindfulness-based relapse prevention [11]. These findings underscore not only the potential wealth of information relating to relapse risk that an individual's location can provide,

but also demonstrate the proven integration of location information into treatment. Furthermore, geolocation data have been specifically identified as being of particular use in both understanding the precipitants to harmful substance use and its effective treatment [12].

Within the substance use literature, geolocation data have historically been used to examine risky locations, such as the influence of neighborhood characteristics on use [13], [14] and individual physical proximity to locations of potential or past harmful use such as bars (either estimated using geofencing or user-defined) [15]–[19]. Several of the applications implemented in these studies enable real-time notifications about locations to their users (e.g., a pop-up message on a smartphone which reads “*You are entering a high-risk zone*”).

On the other hand, affective scientists have focused instead more closely on factors relating to mood. Geolocation data have been used to estimate loneliness and isolation [20], to demonstrate increases in positive affect from seeking out novel environments [21], and to quantify depressive symptoms [22]. Moreover, these data have not only been harnessed to measure mood symptoms, but to also predict their emergence [23].

An integration across these subfields can be in part accomplished by enriching geolocation data with brief, intermittent surveys probing specific information about frequently visited locations. For example, some of the more nuanced facets captured within location are associations with others (or lack thereof, e.g., social isolation), associations with previous drinking behaviors (e.g., whether or not alcohol is present), and associations with affect (i.e., negative versus positive emotions tied to a given location).

Model Evaluation

Data selection, however, is only one component of the successful development of a continuous risk monitoring algorithm. Following algorithm development, it is imperative that these models be rigorously evaluated using performance metrics and eventually tested using independent observations (i.e., using data from individuals which were not used in model development). This workflow in machine learning is what enables researchers to anticipate how well a model could be expected to generalize to new populations and is key when aiming to develop algorithms for real-world healthcare implementation. While standard performance metrics like model accuracy, for example, have been standard reporting practice for years, recent literature has begun to urge researchers to also include assessments of how *fair* a model is [24], [25]. A fair algorithm is one with no preference in performance with respect to inherent or acquired characteristics (e.g., gender, race, socioeconomic status; [26]). In the context of a continuous risk monitoring algorithm for AUD, this would mean that lapse predictions are reasonably accurate and do not favor or disadvantage any particular group.

The motivating factors behind this call to action are clear. In the broader context of health-related data, historical patterns of health care inequities will almost certainly and unavoidably be embedded within data used to train algorithms. These inequities may unintentionally be carried forward in perpetuity by machine learning models if not critically examined. Without examining algorithmic fairness prior to deployment in the real-world, monitoring algorithms run the risk of providing sub-optimal mental health care to individuals who already face disadvantages.

For example, having a limited number of observations within underrepresented groups means that our models will not have as wide a range of individuals to learn from for making predictions of lapse as compared to white, non-Hispanic participants. Performance of these models for racialized minority individuals may therefore be less accurate as a result, particularly without the use of resampling techniques to amend these imbalances [27], [28].

Fairness is also a particular concern specifically in the context of AUD, where the literature has historically been built upon research developed with male, predominantly white, participants. Despite the call to action brought forth by the NIH through their *Guidelines on Inclusion of Women and Minorities in Research*, recent work has highlighted that seminal research in the field on medications for the treatment of AUD have failed to consistently report participant demographics [29]. This lack of reporting makes it difficult to assess how and if this lack of representation is being corrected. By the very nature of its historically limited participant pool, AUD research and its theory have been developed from a particular perspective using a particular group of individuals. This means that the variables that researchers decide are important to measure and input into models, informed by knowledge of AUD theory, will inherently be biased and may favor these groups. Therefore, researchers can also not assume that balanced classes are enough to compensate for biases brought on by the broader societal context. Both of these facts motivate the reasoning behind examining algorithmic fairness in the context of developing continuous risk monitoring systems.

The Current Study

In order for these continuous risk monitoring systems to be implemented in the real-world, these models must both be developed outright and rigorously evaluated on both standard performance metrics and algorithmic fairness. To this end, this study utilized geolocation data collected from smartphones and corresponding self-reported contextual information for frequently visited locations to build a machine learning model to predict next-day alcohol use lapse among individuals with a diagnosis of AUD and a recovery goal of abstinence. Model features were engineered from both raw geolocation data, both context-dependent and context agnostic, and change in these data over the previous 6, 12, 24, 48, 72, and 168 hours.

Here we present characterization of model performance for this prediction model in a validation set. We also evaluated feature importance and model fairness. This study constitutes a preliminary evaluation of a model designed to predict lapse back to alcohol use using minimally burdensome data that has the potential to be integrated within a continuous risk monitoring platform.

Methods

Participants

One hundred and forty six individuals in early-recovery (1-8 weeks of abstinence) for AUD were recruited from the Madison area to take part in a three-month study on how mobile health technology can provide recovery support between 2017 and 2019 (R01 AA024391). Recruitment approaches included social media platforms (e.g., Facebook), television and radio advertisements,

and clinic referrals. Prospective participants completed a phone screen to assess match with eligibility criteria (Table 1). Participants were excluded if they exhibited severe symptoms of paranoia or psychosis (a score ≤ 2.24 on the SCL-90 psychosis scale or a score ≤ 2.82 on the SCL-90 paranoia scale administered at screening).

| Eligibility Criteria |
|--|
| ≥ 18 years of age |
| Ability to read and write in English |
| Diagnosis of moderate AUD (≥ 4 self-reported DSM-5 symptoms) |
| Abstinent from alcohol for 1-8 weeks |
| Willing to use only one smartphone** while on study |

Table 1: Eligibility criteria for study enrollment. **Personal or study-provided.

Procedure

Participants enrolled in a three-month study consisting of five in-person visits, daily surveys, and continuous passive monitoring of geolocation data. Following screening and enrollment visits in which participants consented to participate, learned how to manage location sharing (i.e., turn off location sharing when desired), and reported frequently visited locations, participants completed three follow-up visits one month apart. At each visit, participants were asked questions about frequently visited (>2 times during the course of the previous month) locations. Participants were debriefed at the third and final follow-up visit. Participants were expected to provide continuous geolocation data while on study. Other personal sensing data streams (EMA, cellular communications, sleep quality, and audio check-ins) were collected as part of the parent grant's aims (R01 AA024391).

Geolocation data

To enable collection of geolocation data, participants downloaded either the Moves app or the FollowMee app during the intake visit. Moves was bought-out and subsequently deprecated while the study was ongoing (July 2018) and data collection continued using FollowMee until the end of the study. Both apps continuously tracked location via GPS and WiFi positioning technology.

Data were then processed to filter out duplicated points, fast movement speeds ($>100\text{mph}$), sudden positional jumps, and periods of long duration suggesting sampling error issues (>24 hours with no movement or >2 hours with a positional jump of more than 0.31 miles or 500 meters). Data points were classified as “in transit” when spacing between individual positions suggested a movement speed of greater than 4mph per NIH health guidelines [30].

Contextual information

Contextual information for frequently visited locations (>2 times in the previous month) was obtained during an interview at each follow-up visit (at month 1, 2, and 3; Table 2). Participants were considered to be at a known contextual location if they were within 0.031 miles (50 meters) of a reported frequently visited location.

| Question | Responses |
|---|--|
| Address | |
| Type of place | Work, School, Volunteer, Health care, Home of a friend, Home of a family member, Liquor store, Errands (e.g., grocery store, post office), Coffee shop or cafe, Restaurant, Park, Bar, Gym or fitness center, AA or recovery meeting, Religious location (e.g., church, mosque, temple), Other |
| Have you drank alcohol here before? | No, Yes |
| Is alcohol available here? | No, Yes |
| How would you describe your experiences here? | Pleasant, Unpleasant, Mixed, Neutral |
| Does being at this location put you at any risk to begin drinking? | No risk, Low risk, Medium risk, High risk |
| Did the participant identify this place as a risky location they are trying to avoid now that they are sober? | No, Yes |

Table 2: Location information collected from frequently visited locations.

Participant characteristics

Participants completed a baseline measure of demographics and other constructs relevant to lapse at the screening visit, which was used for fairness assessments (Table 3).

| Variable | Measure |
|--------------|---|
| Demographics | Age |
| | Sex |
| | Race |
| | Ethnicity |
| | Employment |
| | Income |
| | Marital Status |
| Alcohol | Alcohol Use History |
| | DSM-5 Checklist for AUD |
| | Young Adult Alcohol Problems Test |
| | WHO-The Alcohol, Smoking and Substance Involvement Screening Test |

Table 3: Demographic and relevant alcohol use history variables sampled at screening visit.

Lapses

Alcohol lapses were used as the outcome variable in this study and were used to provide labels for model training, for testing model performance, and for testing issues of algorithmic fairness across our predefined subgroups. Future lapse occurrence (here conceptualized as next-day lapse) was be predicted in 24-hour windows, beginning at 4:00am on a participant’s second day of participation to ensure one full day of data collection for the first window, and at every subsequent day on study thereafter. *Lapse* and *no lapse* occurrences were identified from the daily survey question, “*Have you drank any alcohol that you have not yet reported?*”.

Participants who responded *yes* to this question were then asked to report the date and hour of the start and the end of the drinking episode. In this case, the prediction window was labeled *lapse*.

Prediction windows were labeled *no lapse* if no alcohol use was reported within that window.

Feature engineering

Feature engineering is the process of creating variables (or “*features*”) from unprocessed data and was used to transform raw data from geolocation data collected the prior day.

Features from geolocation data were generated that utilized both contextual information collected from monthly surveys (e.g., location valence, perceived riskiness) as well as features that were independent of further individual input (e.g., location variance, time spent out of the home in the evening). All features were calculated both as raw and change features based on previous geolocation data (i.e., change from past 6, 12, 24, 48, 72, and 168 hour periods) in order to capture individual variation.

Imputation of missing data and removal of zero-variance features are additional general processing steps that will also be undertaken during feature engineering.

Algorithm development & performance

Several configurations of the XGBoost machine learning algorithm were considered which varied across a relevant and appropriate range of model-specific hyperparameters (mtry, tree depth, learning rate) as well as resampling techniques (up-sampling of the positive class, lapse, and down-sampling of the negative class, no lapse).

Models were trained and assessed using participant-grouped, nested k -fold cross-validation. Grouped cross-validation assigns all data from a participant as either held-in or held-out to avoid bias introduced when predicting a participant’s data from their own data. Nested cross-validation uses two nested loops for dividing and holding out folds: an outer loop, where held-out folds serve as test sets for model evaluation; and inner loops, where held-out folds serve

as validation sets for model selection. Importantly, these sets are independent, maintaining separation between data used to train the models, select the best models, and evaluate those best models. Therefore, nested cross-validation removes optimization bias from the evaluation of model performance in the test sets and can yield lower variance performance estimates than single test set approaches [31].

The primary performance metric for model selection and evaluation of the validation set was area under the Receiver Operating Characteristic Curve (auROC) [32]. auROC indexes the probability that the model will predict a higher score for a randomly selected positive case (lapse) relative to a randomly selected negative case (no lapse). This metric was selected because it 1) combines sensitivity and specificity, which are both important characteristics for clinical implementation; 2) is an aggregate metric across all decision thresholds, which is important because optimal decision thresholds may differ across settings and goals; and 3) is unaffected by class imbalance, which is important for comparing models with differing prediction window widths and levels of class imbalance. The best model configuration was selected using median auROC across all validation sets. Several secondary performance metrics including sensitivity, specificity, balanced accuracy, positive predictive value (PPV), and negative predictive value (NPV) will also be assessed.

SHAP (SHapley Additive exPlanations) values were computed as interpretability metrics to identify the relative importance of different features in each final algorithm. SHAP values measure the unique contribution of features in an algorithm's predictions [33]. SHAP values possess several useful properties including: *Additivity* (SHAP values for each feature can be

computed independently and summed); *Efficiency* (the sum of SHAP values across features must add up to the difference between predicted and observed outcomes for each observation); *Symmetry* (SHAP values for two features should be equal if the two features contribute equally to all possible coalitions); and *Dummy* (a feature that does not change the predicted value in any coalition will have a SHAP value of 0). Highly important features represent relevant, actionable potential antecedents to lapse (and therefore points of intervention) that will be relevant in the future development of a continuous risk monitoring system. However, these are descriptive analyses because standard errors or other indices of uncertainty for importance scores are not available for SHAP values.

Finally, a Bayesian hierarchical generalized linear model was used to estimate the posterior probability distributions and 95% Bayesian credible intervals (CIs) for auROC.

Algorithmic fairness

Subgroups were defined on the basis of personal individual characteristics (in the machine learning fairness literature, “sensitive attributes”) that are specifically associated with treatment disparities in AUD.

A Bayesian hierarchical generalized linear model was used to estimate the posterior probability distributions and 95% Bayesian CIs for auROC across four subgroups of participants: white versus non-white, younger than 55 versus equal to and older than 55, below or above the federal poverty line (citation), and sex (male versus female).

Results

Demographics

A total of 192 individuals were eligible to participate in the study, of which 191 consented to participate and 169 enrolled in the study. Fifteen participants were excluded prior to the first monthly follow-up visit. One participant was excluded for not maintaining a recovery goal of abstinence during their time on study. Two participants were excluded due to evidence of low compliance and careless responding. A further five individuals were excluded due to poor geolocation data quality as a result of insufficient data (primarily resulting from software incompatibility), resulting in a final sample size of 146.

The average age of the final sample was 40.9 years (SD = 12 years, range = 21-72 years). There was an approximately equal number of men (n = 74, 50.7%) and women (n = 72, 49.3%). The majority of the sample was White/Caucasian (n = 127, 86.99%) and non-Hispanic (n = 142, n = 97%). The mean income of participants was \$34,408 (SD = \$32,259, range = \$0-\$200,000). On average, participants self-reported a mean number of 8.9 DSM-V symptoms of AUD (range = 4-11).

Model Evaluation

We selected and evaluated the best performing XGBoost model from our validation set. This may result in a slight optimism bias in our model performance, though we believe this is largely offset through our use of 10 x 30 cross-validation (which averages model performance across 300 folds). However, evaluating the validation set was important to do because model development is still in progress, and as such it would not have been appropriate to examine independent test set performance at this stage.

need to pull stats!

The median auROC across the 300 folds achieved fair performance (Mdn = 0.721, IQR = XX, range = XX-XX). Figure 1 displays a histogram of model performance distribution across all folds. A receiver operating characteristic curve is displayed in Figure 2, representing aggregate predicted lapse probabilities across all 300 folds of the validation set.

We also used the posterior probability distributions for the auROCs for our best performing validation set model to formally characterize model performance. The median auROC was XX (95% CI [XX-XX]; Figure 3).

Finally, we performed model calibration in order to improve our trust in model predictions. Results of model calibration are displayed in Figure 4, showing that this model *over* predicts lapse probability even after calibrating the model. In other words, our model is more likely to predict that an individual will lapse than the true rate of lapse in our sample.

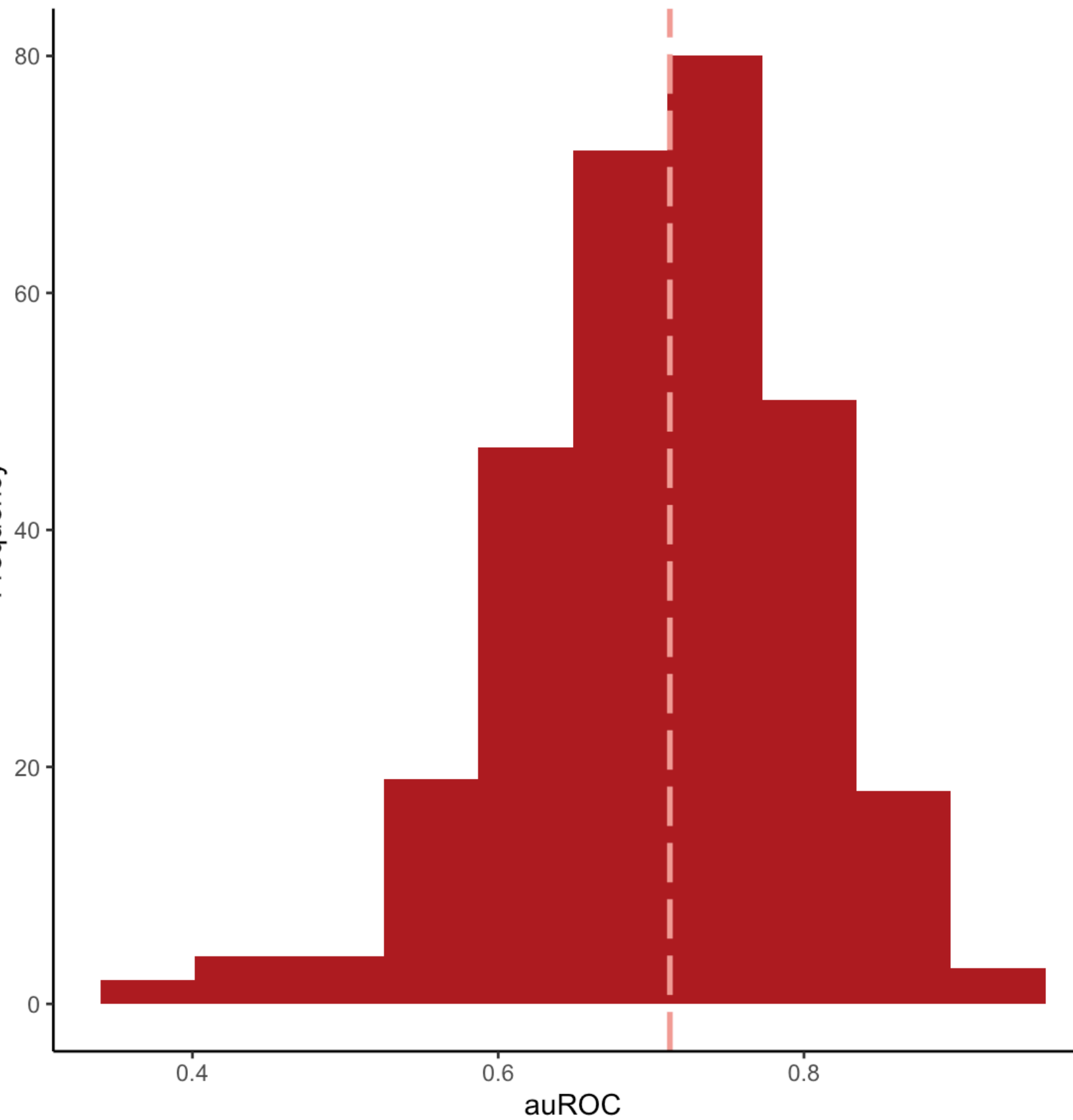


Figure 1: auROC performance across all inner cross-validation folds.

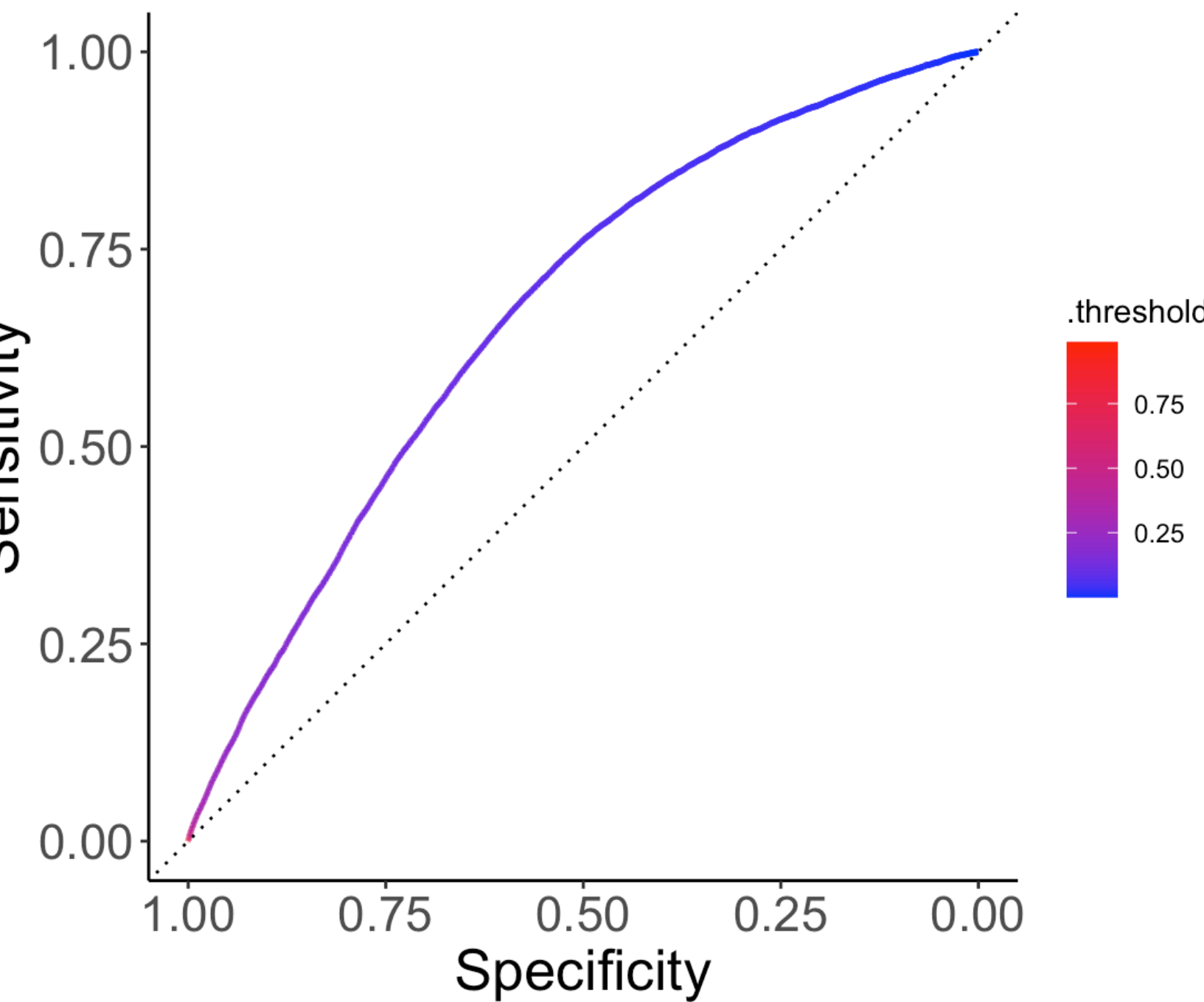


Figure 2: auROC for overall validation set performance across thresholds.

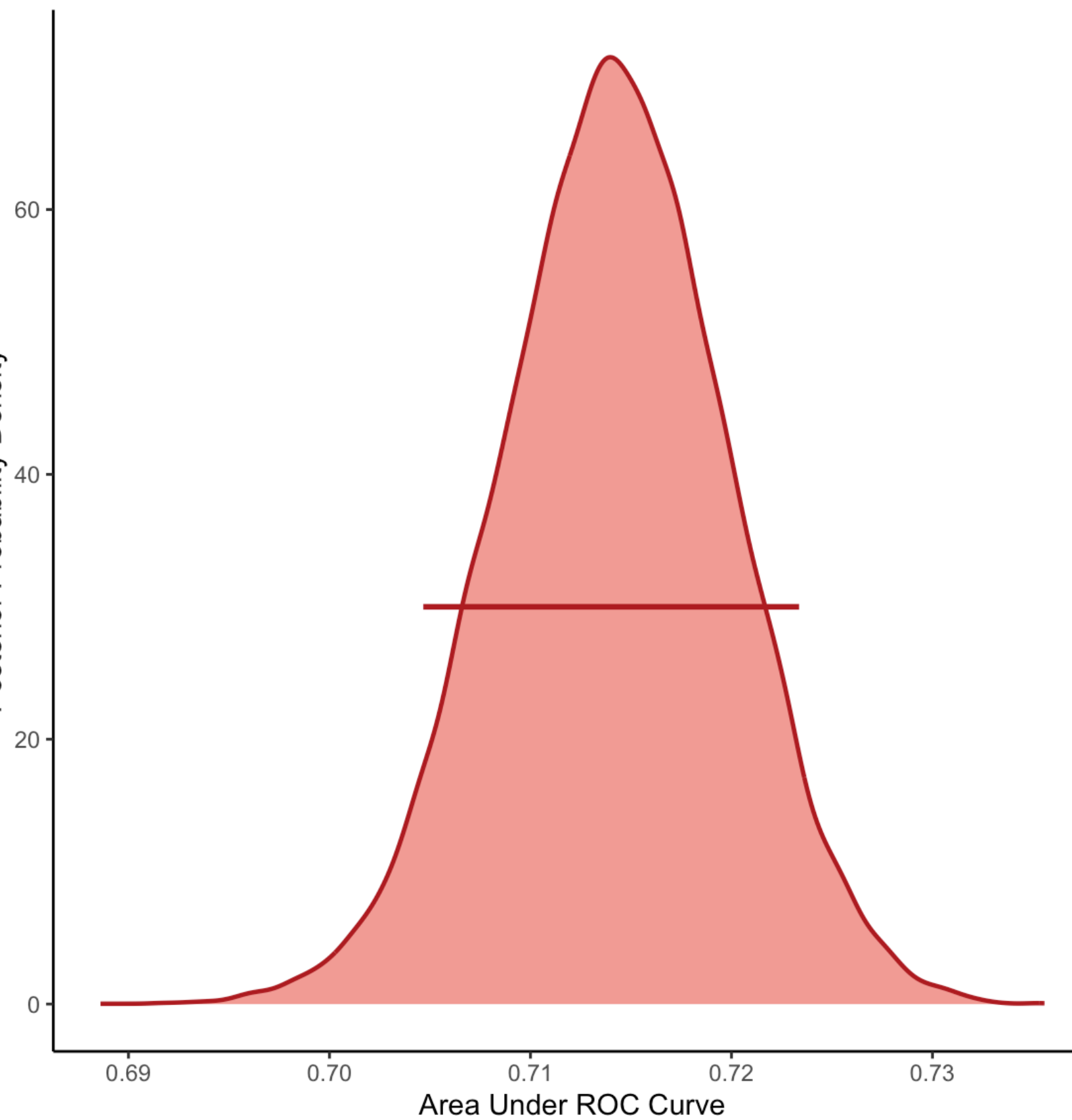


Figure 3: Posterior distribution of model performance with 95% credible interval.

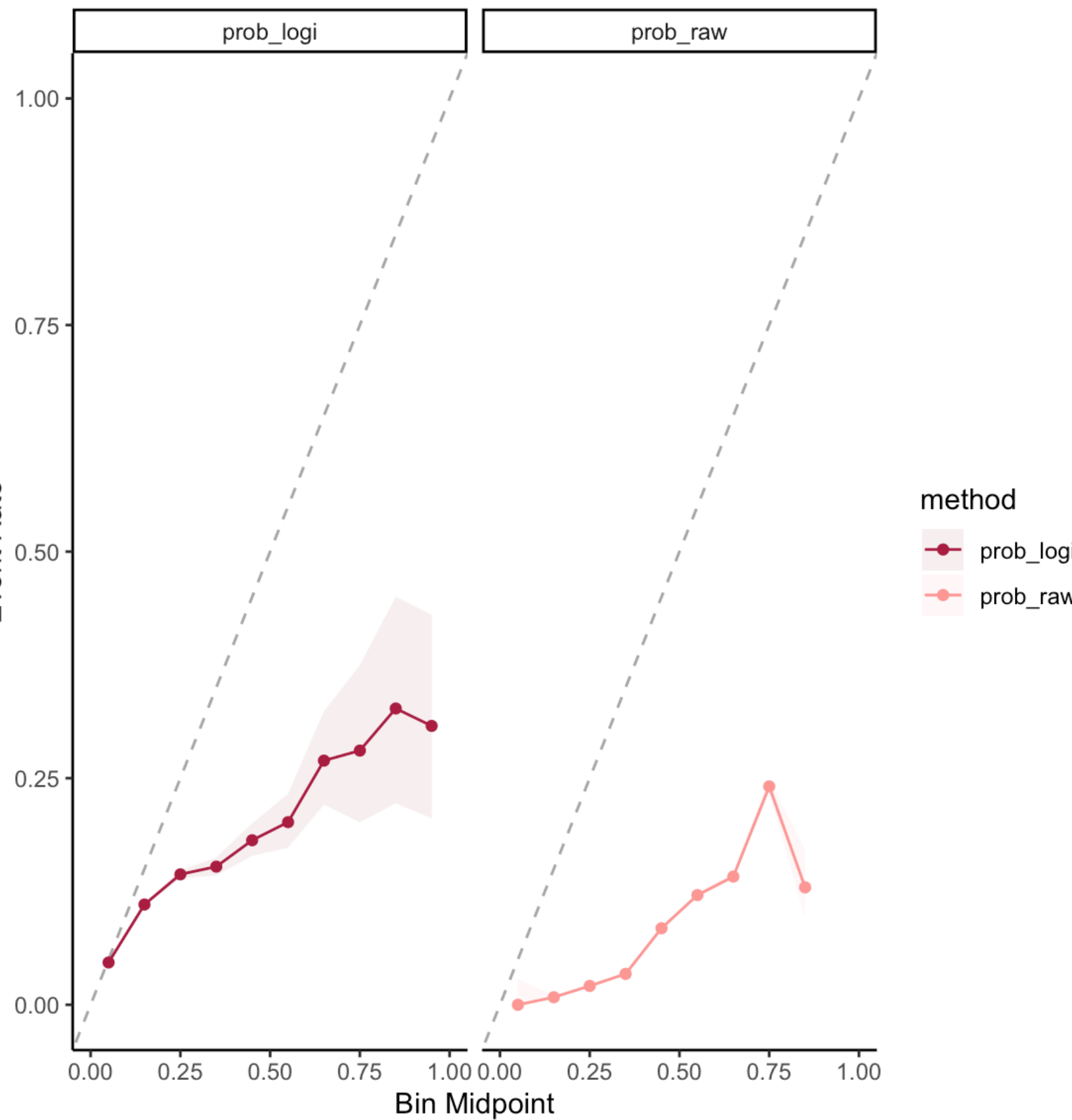


Figure 4: Comparison between raw (uncalibrated) and logistic (calibrated) probabilities.

Feature Importance

Global importance (mean absolute Shapley values) for feature categories is shown in Figure 5.

Three features were identified as being particularly important in contributing to model predictions: time spent at risky locations, time spent at different types of location, and time spent at locations with varying levels of alcohol availability. Other features, both context-supplemented and without, did not appear to be strong global contributors to model predictions.

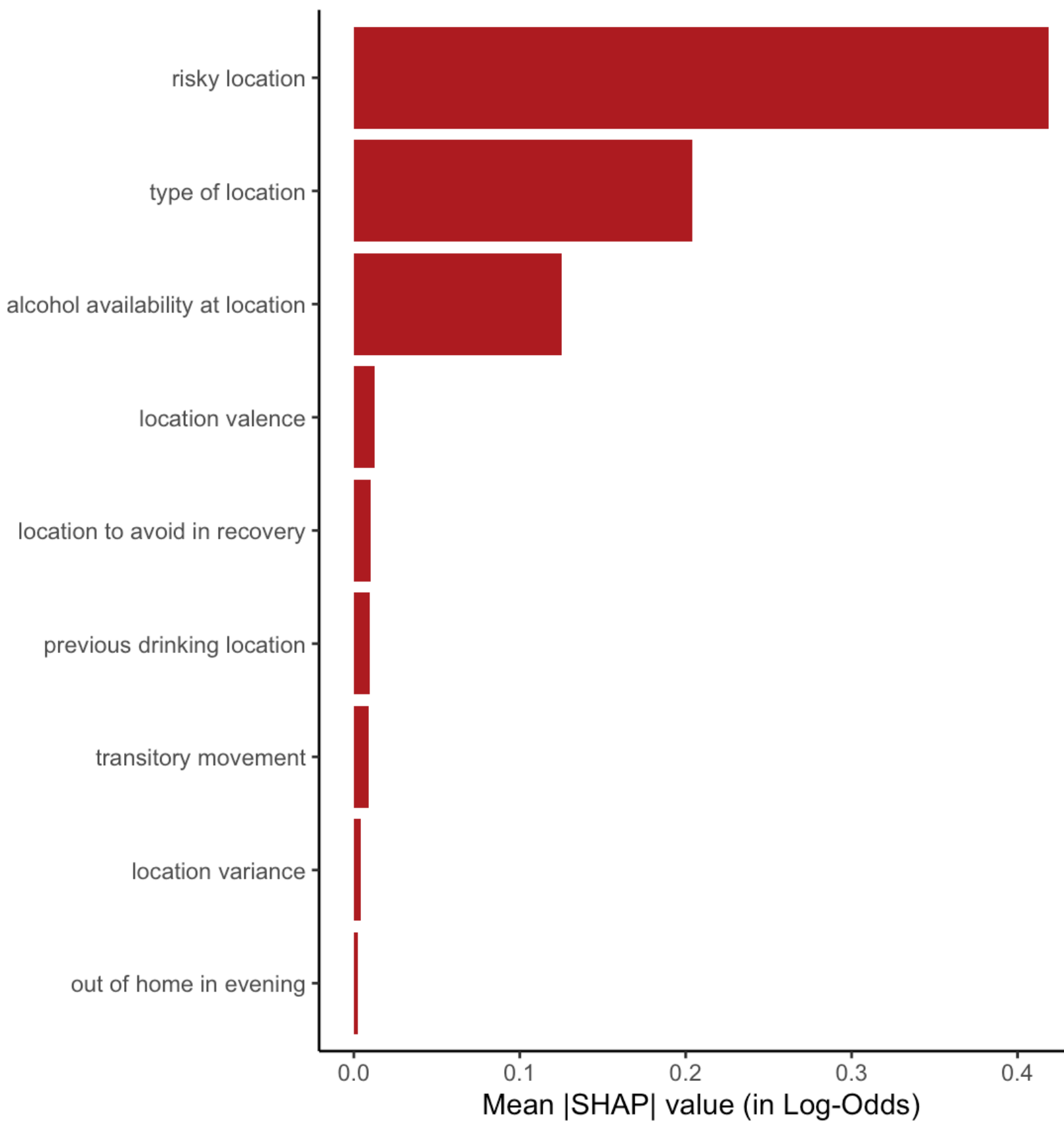


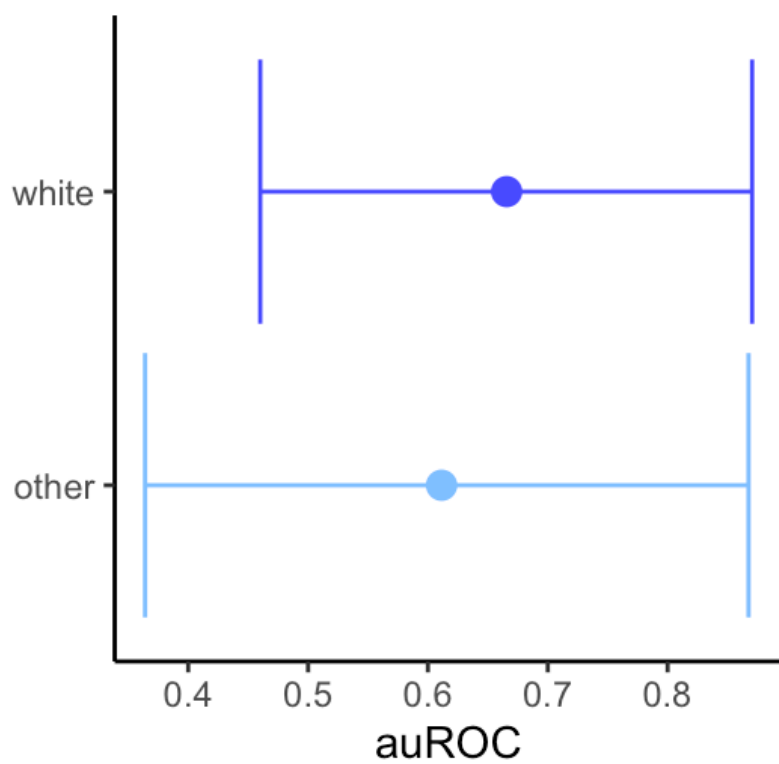
Figure 5: Grouped (mean absolute values) SHAP values displaying relative feature importance.

SINA plot?

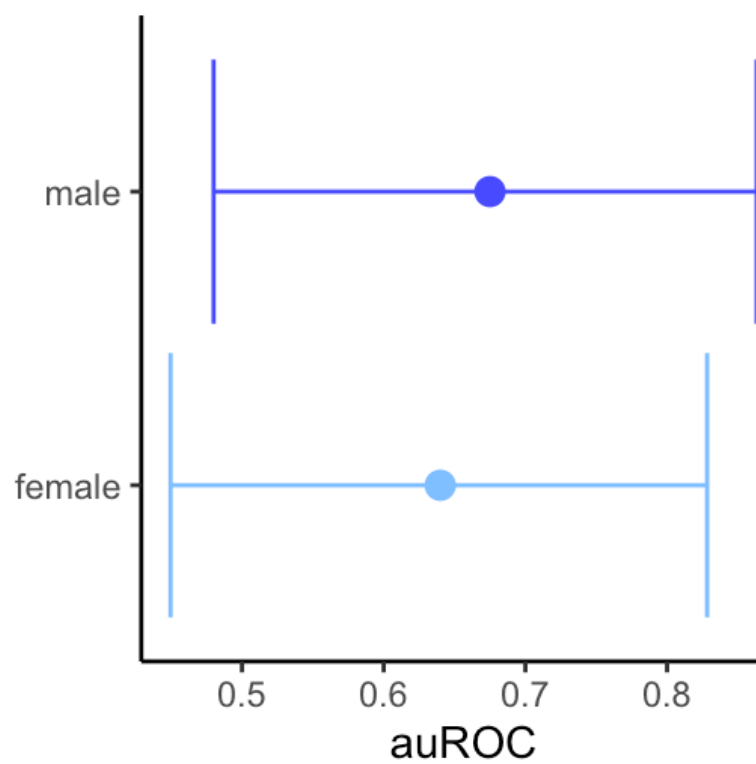
Algorithmic Fairness

Figure 6 shows the differences in model performance across race (white = 127, non-white = 19), sex (male N = 74, female N = 72), age (younger than 55 = 126, older than or equal to 55 = 20), and income (below federal poverty line = 48, above federal poverty line = 98). No group comparisons were significant (probability < .95) across models, suggesting comparable model performance across subgroups.

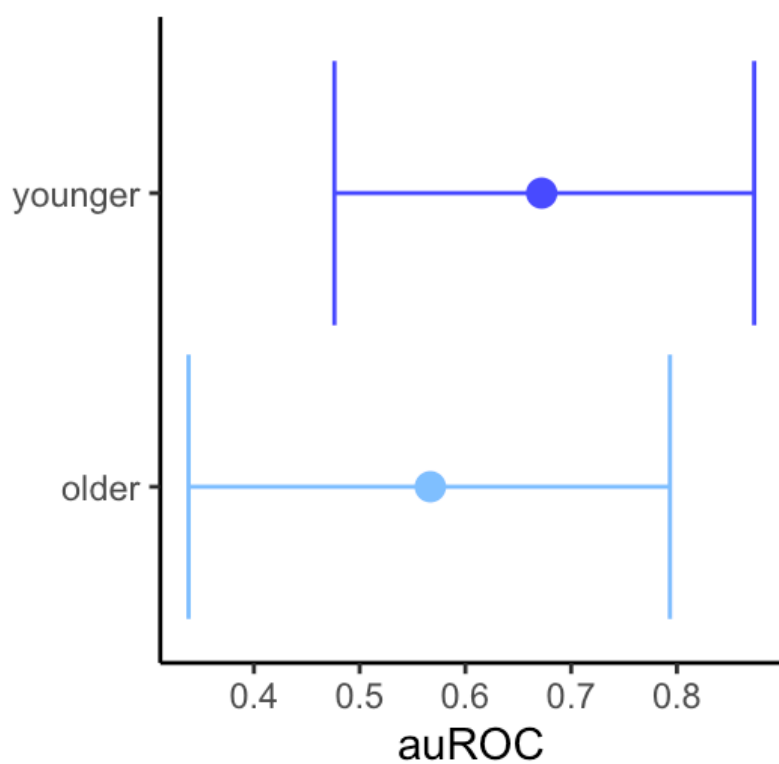
Race/Ethnicity



Sex



Age



Income

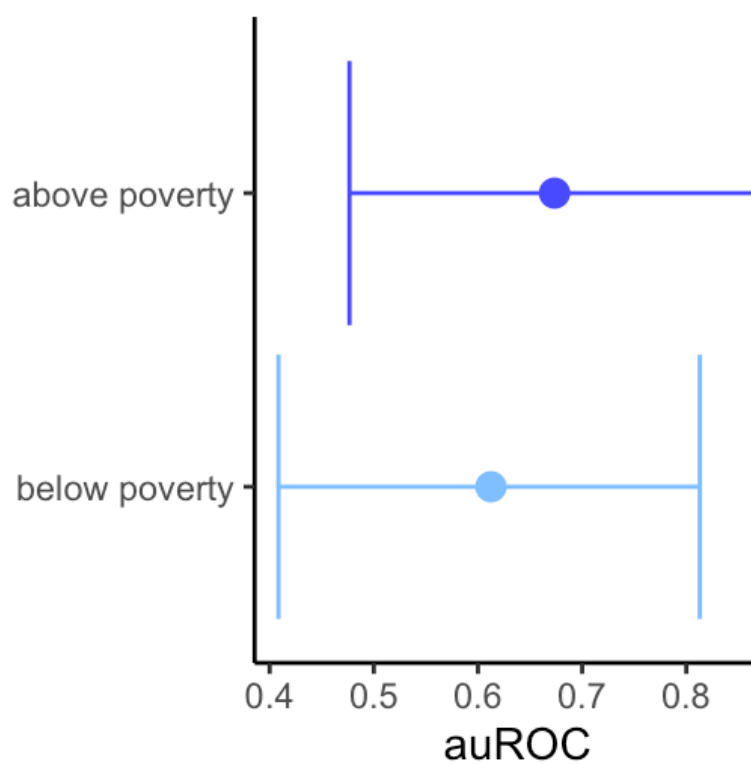


Figure 6: 95% credible intervals across posterior subgroups at differential levels of privilege.

Discussion

Model Performance

- Our day-level model of lapse prediction using geolocation data performs adequately well.

Models which perform at around the .7 threshold are considered to have “fair” performance.

This suggests that, while there is still substantial room for improvement in model performance, geolocation data can predict future alcohol lapse in the next day with fair sensitivity and specificity.

- Calibration is more sensitive than true lapse rate in sample, suggesting that it overpredicts occurrence of lapses
- This is not necessarily an issue if we are trying to quantify relative risk to individuals using a risk monitoring system
- Top performing Shapley values were time spent at risky locations, time spent at different types of location, and time spent at locations with varying levels of alcohol availability.
- These features were all generated utilizing additional context supplied by participants after a given location was identified as frequently visited ($> 2x$ in the previous month).
- It should be noted that these features have potential to be generated without user feedback (i.e., using consumer and other publicly available data to identify establishments that sell alcohol, etc.)
- This could potentially reduce burden further by not requiring individual input
- However, self-classifying locations as risky might be encoding nuance that could not be feasibly obtained using public data. For example, a location might be labeled as risky from user input because it is a person-specific triggering location (e.g., scene of a traumatic event).

- Interestingly, location valence (i.e., the emotion tied to a given location) is the fourth-highest Shapley value, yet appears to be minimally contributing to model predictions. This may be because participants were asked retrospectively about these locations at one month follow-up visits, and so our measures of emotional quality of a location may be too distal to be meaningful (particularly when compared to daily EMA).
- Location to avoid in recovery and previous drinking location – poor insight so low predictive ability?
- Transitory movement, location variance, time spent out of home in the evening – lacking individual input from participants, not easily tied to meaningful lapse risk factors

Model Fairness

- Performance appears to be comparable across subgroups, but may be due to only fair performance of the model
- We could be more sure of the conclusions from this comparison if we had a model with excellent performance
- Fairness analyses will need to be repeated on the final sample

Future directions

- Baseline model?
- Add in more affective features
- Add in risk-terrain modeling features
- Add in other important features that could contribute to movement patterns like day of the week and weather

- Test final model
- Further calibration
- Break down three top performing features into their subcomponents (i.e., high risk locations, medium risk locations, and low risk locations; yes there is alcohol available here or no there is not alcohol available here) to obtain a more nuanced understanding of model performance.

Conclusion

This study demonstrates that it is feasible to predict lapse with a fair level of accuracy using geolocation data, suggesting that geolocation data is a viable supplement for risk prediction monitoring systems. Moreover, our model demonstrates similar performance across vulnerable subgroups. Moving forward, additional risk-relevant features will be added to the model in an effort to improve prediction and the final model will be evaluated.

References

Bibliography

- [1] “Highlights for the 2022 National Survey on Drug Use and Health”.
- [2] J. A. Tucker, S. D. Chandler, and K. Witkiewitz, “Epidemiology of Recovery From Alcohol Use Disorder”, *Alcohol Research : Current Reviews*, vol. 40, no. 3, p. 2, Nov. 2020, doi: 10.35946/arcv.v40.3.02.
- [3] “The Science of Drug Use and Addiction: The Basics | NIDA Archives”.

- [4] T. H. Brandon, J. I. Vidrine, and E. B. Litvin, “Relapse and Relapse Prevention”, *Annual Review of Clinical Psychology*, vol. 3, no. Volume3, 2007, pp. 257–284, Apr. 2007, doi: 10.1146/annurev.clinpsy.3.022806.091455.
- [5] K. Witkiewitz and G. A. Marlatt, “Relapse Prevention for Alcohol and Drug Problems: That Was Zen, This Is Tao”, *The American Psychologist*, vol. 59, no. 4, pp. 224–235, 2004, doi: 10.1037/0003-066X.59.4.224.
- [6] D. C. Mohr, M. Zhang, and S. M. Schueller, “Personal Sensing: Understanding Mental Health Using Ubiquitous Sensors and Machine Learning”, *Annual Review of Clinical Psychology*, vol. 13, no. 1, pp. 23–47, May 2017, doi: 10.1146/annurev-clinpsy-032816-044949.
- [7] P. A. Areàn, K. Hoa Ly, and G. Andersson, “Mobile Technology for Mental Health Assessment”, *Dialogues in Clinical Neuroscience*, vol. 18, no. 2, pp. 163–169, Jun. 2016.
- [8] P. H. Janak and N. Chaudhri, “The Potent Effect of Environmental Context on Relapse to Alcohol-Seeking After Extinction”, *The Open Addiction Journal*, vol. 3, pp. 76–87, Jan. 2010, doi: 10.2174/1874941001003010076.
- [9] M. A. Walton, F. C. Blow, C. R. Bingham, and S. T. Chermack, “Individual and Social/Environmental Predictors of Alcohol and Drug Use 2 Years Following Substance Abuse Treatment”, *Addictive Behaviors*, vol. 28, no. 4, pp. 627–642, Jun. 2003, doi: 10.1016/s0306-4603(01)00284-2.
- [10] M. A. Walton, T. M. Reischl, and C. S. Ramanathan, “Social Settings and Addiction Relapse”, *Journal of Substance Abuse*, vol. 7, no. 2, pp. 223–233, 1995, doi: 10.1016/0899-3289(95)90006-3.

- [11] M. R. LeCocq, P. A. Randall, J. Besheer, and N. Chaudhri, “Considering Drug-Associated Contexts in Substance Use Disorders and Treatment Development”, *Neurotherapeutics: The Journal of the American Society for Experimental NeuroTherapeutics*, vol. 17, no. 1, pp. 43–54, Jan. 2020, doi: 10.1007/s13311-019-00824-2.
- [12] G. J. Stahler, J. Mennis, and D. A. Baron, “Geospatial Technology and the "Exposome": New Perspectives on Addiction”, *American Journal of Public Health*, vol. 103, no. 8, pp. 1354–1356, Aug. 2013, doi: 10.2105/AJPH.2013.301306.
- [13] D. H. Epstein *et al.*, “Real-Time Tracking of Neighborhood Surroundings and Mood in Urban Drug Misusers: Application of a New Method to Study Behavior in Its Geographical Context”, *Drug and Alcohol Dependence*, vol. 134, pp. 22–29, Jan. 2014, doi: 10.1016/j.drugalcdep.2013.09.007.
- [14] M.-P. Kwan, J. Wang, M. Tyburski, D. H. Epstein, W. J. Kowalczyk, and K. L. Preston, “Uncertainties in the Geographic Context of Health Behaviors: A Study of Substance Users' Exposure to Psychosocial Stress Using GPS Data”, *International Journal of Geographical Information Science*, vol. 33, no. 6, pp. 1176–1195, Jun. 2019, doi: 10.1080/13658816.2018.1503276.
- [15] S. Attwood, H. Parke, J. Larsen, and K. L. Morton, “Using a Mobile Health Application to Reduce Alcohol Consumption: A Mixed-Methods Evaluation of the Drinkaware Track & Calculate Units Application”, *BMC public health*, vol. 17, no. 1, p. 394, May 2017, doi: 10.1186/s12889-017-4358-9.
- [16] S. Carreiro, M. Taylor, S. Shrestha, M. Reinhardt, N. Gilbertson, and P. Indic, “Realize, Analyze, Engage (RAE): A Digital Tool to Support Recovery from Substance Use

Disorder”, *Journal of Psychiatry and Brain Science*, vol. 6, p. e210002, 2021, doi: 10.20900/jpbs.20210002.

- [17] V. M. Gonzalez and P. L. Dulin, “Comparison of a Smartphone App for Alcohol Use Disorders with an Internet-based Intervention plus Bibliotherapy: A Pilot Study”, *Journal of Consulting and Clinical Psychology*, vol. 83, no. 2, pp. 335–345, Apr. 2015, doi: 10.1037/a0038620.
- [18] D. H. Gustafson *et al.*, “A Smartphone Application to Support Recovery From Alcoholism: A Randomized Clinical Trial”, *JAMA Psychiatry*, vol. 71, no. 5, p. 566, May 2014, doi: 10.1001/jamapsychiatry.2013.4642.
- [19] F. Naughton *et al.*, “A Context-Sensing Mobile Phone App (Q Sense) for Smoking Cessation: A Mixed-Methods Study”, *JMIR mHealth and uHealth*, vol. 4, no. 3, p. e106, Sep. 2016, doi: 10.2196/mhealth.5787.
- [20] A. Doryab *et al.*, “Identifying Behavioral Phenotypes of Loneliness and Social Isolation with Passive Sensing: Statistical Analysis, Data Mining and Machine Learning of Smartphone and Fitbit Data”, *JMIR mHealth and uHealth*, vol. 7, no. 7, p. e13209, Jul. 2019, doi: 10.2196/13209.
- [21] A. S. Heller *et al.*, “Association between Real-World Experiential Diversity and Positive Affect Relates to Hippocampal-Striatal Functional Connectivity”, *Nature Neuroscience*, vol. 23, no. 7, pp. 800–804, Jul. 2020, doi: 10.1038/s41593-020-0636-4.
- [22] I. M. Raugh *et al.*, “Geolocation as a Digital Phenotyping Measure of Negative Symptoms and Functional Outcome”, *Schizophrenia Bulletin*, vol. 46, no. 6, pp. 1596–1607, Dec. 2020, doi: 10.1093/schbul/sbaa121.

- [23] J. Shin and S. M. Bae, “A Systematic Review of Location Data for Depression Prediction”, *International Journal of Environmental Research and Public Health*, vol. 20, no. 11, p. 5984, May 2023, doi: 10.3390/ijerph20115984.
- [24] A. Rajkomar, M. Hardt, M. D. Howell, G. Corrado, and M. H. Chin, “Ensuring Fairness in Machine Learning to Advance Health Equity”, *Annals of Internal Medicine*, vol. 169, no. 12, pp. 866–872, Dec. 2018, doi: 10.7326/M18-1990.
- [25] J. Wawira Gichoya, L. G. McCoy, L. A. Celi, and M. Ghassemi, “Equity in Essence: A Call for Operationalising Fairness in Machine Learning for Healthcare”, *BMJ health & care informatics*, vol. 28, no. 1, p. e100289, Apr. 2021, doi: 10.1136/bmjhci-2020-100289.
- [26] X. Wang, Y. Zhang, and R. Zhu, “A Brief Review on Algorithmic Fairness”, *Management System Engineering*, vol. 1, no. 1, p. 7, Nov. 2022, doi: 10.1007/s44176-022-00006-z.
- [27] N. Japkowicz, “The Class Imbalance Problem: Significance and Strategies”, in *Proc. of the Int'l Conf. on Artificial Intelligence*, 2000, pp. 111–117.
- [28] A. Wang, V. V. Ramaswamy, and O. Russakovsky, “Towards Intersectionality in Machine Learning: Including More Identities, Handling Underrepresentation, and Performing Evaluation”, in *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, in FAccT '22. New York, NY, USA: Association for Computing Machinery, Jun. 2022, pp. 336–349. doi: 10.1145/3531146.3533101.
- [29] M. R. Schick, N. S. Spillane, and K. L. Hostetler, “A Call to Action: A Systematic Review Examining the Failure to Include Females and Members of Minoritized Racial/Ethnic Groups in Clinical Trials of Pharmacological Treatments for Alcohol Use Disorder”,

Alcoholism: Clinical and Experimental Research, vol. 44, no. 10, pp. 1933–1951, 2020, doi: 10.1111/acer.14440.

[30] “Physical Activity Guidelines for Americans, 2nd Edition”.

[31] P. Jonathan, W. J. Krzanowski, and W. V. McCarthy, “On the Use of Cross-Validation to Assess Performance in Multivariate Prediction”, *Statistics and Computing*, vol. 10, no. 3, pp. 209–229, Jul. 2000, doi: 10.1023/A:1008987426876.

[32] M. Kuhn and K. Johnson, *Applied Predictive Modeling*, 1st ed. 2013, Corr. 2nd printing 2018 edition. New York: Springer, 2018. doi: 10.1007/978-1-4614-6849-3.

[33] S. M. Lundberg and S.-I. Lee, “A Unified Approach to Interpreting Model Predictions”, in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, in NIPS'17. Red Hook, NY, USA: Curran Associates Inc., Dec. 2017, pp. 4768–4777.