

Table 1: Best Balanced Accuracy for Each Statistical Algorithm across 10 Folds during Model Selection

algorithm	feature set	feature type(s)	resample	balanced accuracy
random forest	passive	perc, raw	down	0.61
glmnet	active	raw	up	0.60
knn	passive	raw	down	0.58

Table 2: Confusion Matrix for Best Model (passive random forest)

Prediction	Truth	
	no	yes
no	143196	5386
yes	52314	4904

## Results

### Participant Characteristics

### Best Model Performance

We optimized each statistical algorithm by tuning hyperparameter values and fitting models across several feature set combinations (i.e., active or passive, and type of feature engineering). Each model configuration was fit using a grouped 1x10 resampling method. Table 1 shows the best performing model (i.e., highest balanced accuracy) for each statistical algorithm. Figure 1 shows the model’s performance in each held out fold.

Our top performing model, with the highest balanced accuracy, was a random forest statistical algorithm using passive features. To reduce the effects of optimization bias on our model evaluation of predictive performance, we refit the top performing model 100 times (grouped 10x10 resampling). We then averaged across performance estimates to get an estimate with low variance. This method gave us a balanced accuracy estimate of .60. Table 2 shows a confusion matrix where we can see how well the model predicts with negative cases (i.e., no lapse) compared to positive cases (i.e., lapses). Table 3 characterizes our best model over several metrics appropriate for classification.

Since we did not have an independent held out test set we were not able to completely remove optimization bias. So, we performed a model comparison to assess our model’s performance compared to a null model (i.e., intercept only) with a balanced accuracy of .50. A Bayesian correlated t-test revealed a posterior probability that the balanced accuracy of our model was above the Region of Practical Equivalence (ROPE) is .999 (Figure 2). This suggests there is a meaningful difference between our model and a null model with no signal.

In the appendix we show that our best performing model has variation in predictions for each individual participant. Figure A1 contains predictions that are predicted probabilities of a lapse. Each observation was held out 10 times and the figure shows the averaged probability across these 10 predictions. Actual lapses, are depicted in red.

Table 3: Classification Performance Metrics for Best Model (random forest passive) across 100 Folds

metric	estimate
balanced accuracy	0.60
accuracy	0.72
sens	0.46
spec	0.73
ppv	0.09
npv	0.96
area under the ROC curve	0.64

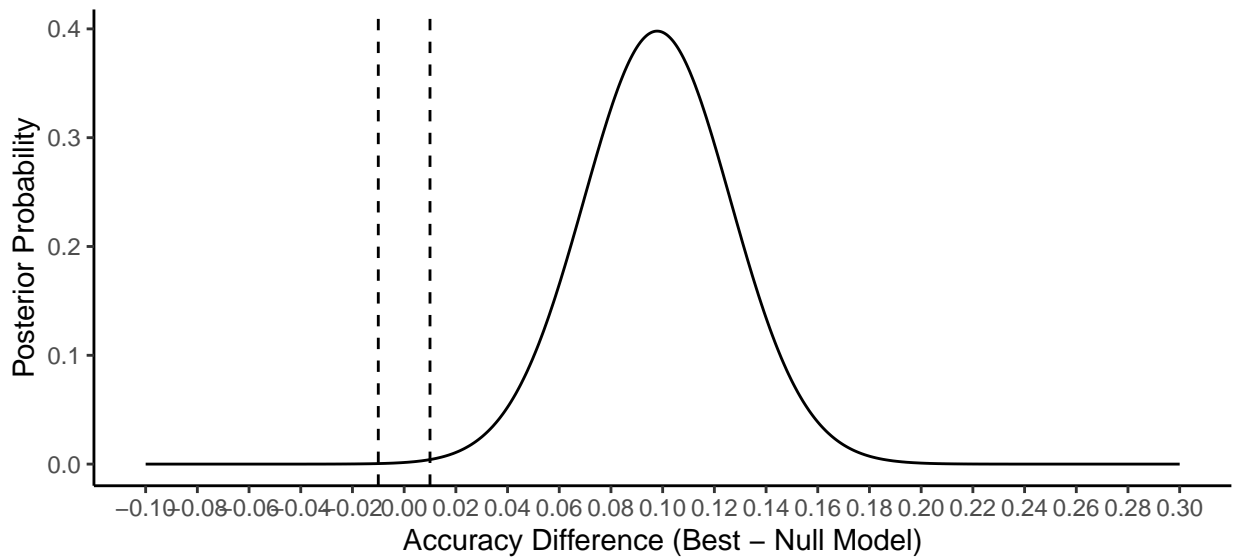


Figure 1. Model comparison of best model (passive random forest) and null model (intercept only).

## Model Comparison between Active and Passive

Our best performing active model using all the features available was a glmnet algorithm with a balanced accuracy of XX. Table 4 characterizes our best active model using various performance metrics. Figure 3 shows the balanced accuracy estimates for passive and active models over 100 fits.

A Bayesian correlated t-test revealed there were no differences between the best active and best passive models. The posterior probability that was between the ROPE was .19 (Figure 4). This suggests we can predict lapses just as well with only passive features compared to models with active features.

Table 4: Classification Performance Metrics for Best Active Model (glmnet) across 100 Folds

metric	estimate
balanced accuracy	0.59
accuracy	0.73
sens	0.44
spec	0.74
ppv	0.09
npv	0.96
area under the ROC curve	0.63

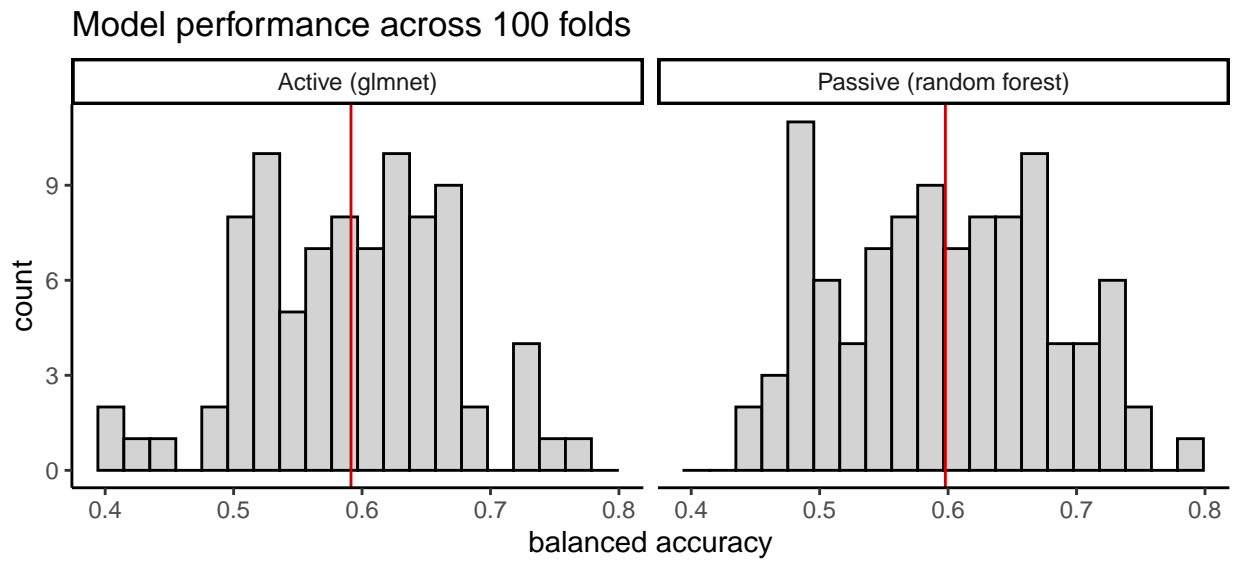


Figure 2. Histogram of model fits across 100 folds. (FIX: 14 folds still running so only 86 splits in active)

```
## Warning in cv_full - cv_compact: longer object length is not a multiple of
## shorter object length
```

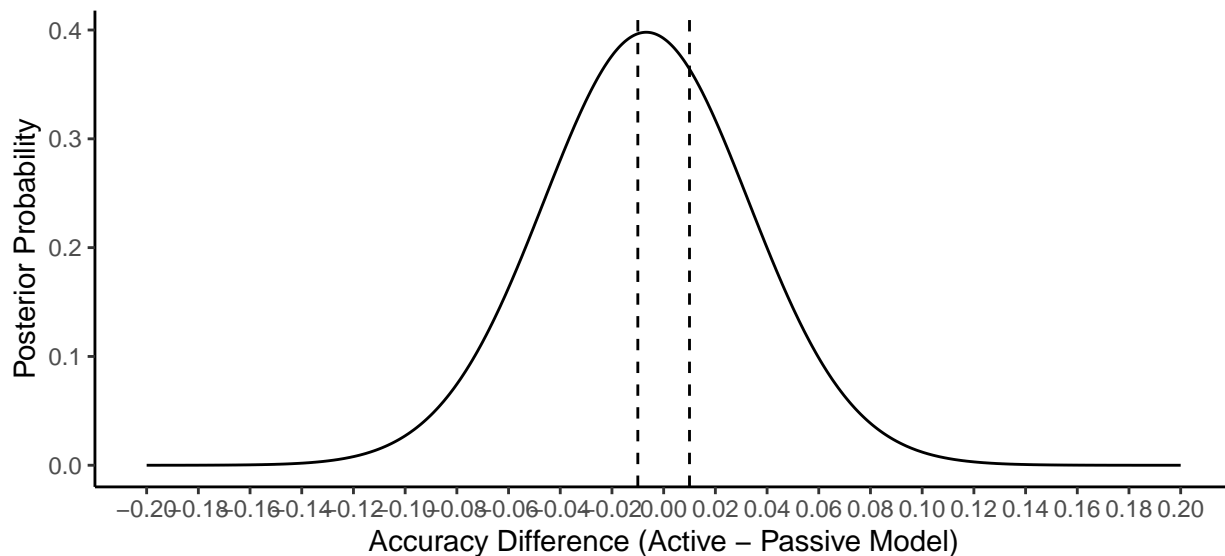


Figure 3. Model comparison of the best passive (random forest) and best active (glmnet) models.

## Top Features

We conducted a permutation test of feature importance to see which features contributed most to our top performing model. Figure 4 shows which features when removed from the model had the greatest effect on performance. Figure 5 shows the top features for the best active model.

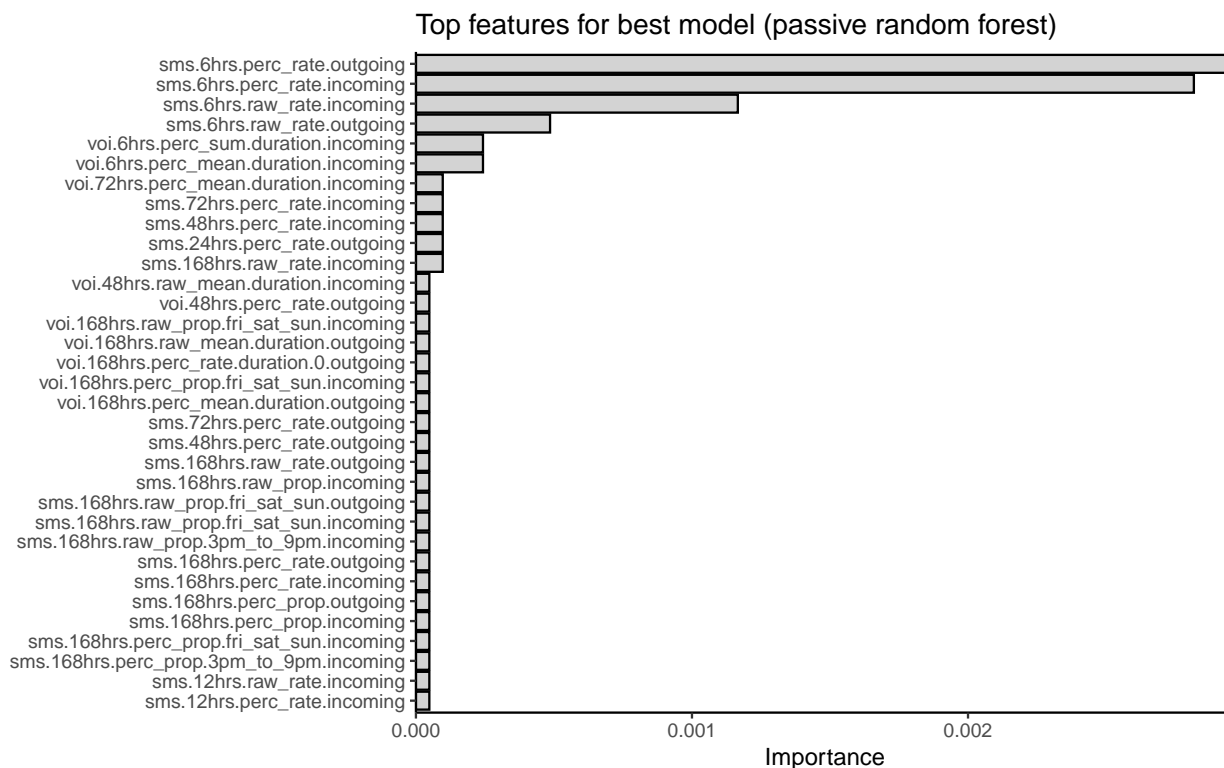


Figure 4. Feature importance scores for best performing model (Passive Random Forest).

Top features for best active model (glmnet)

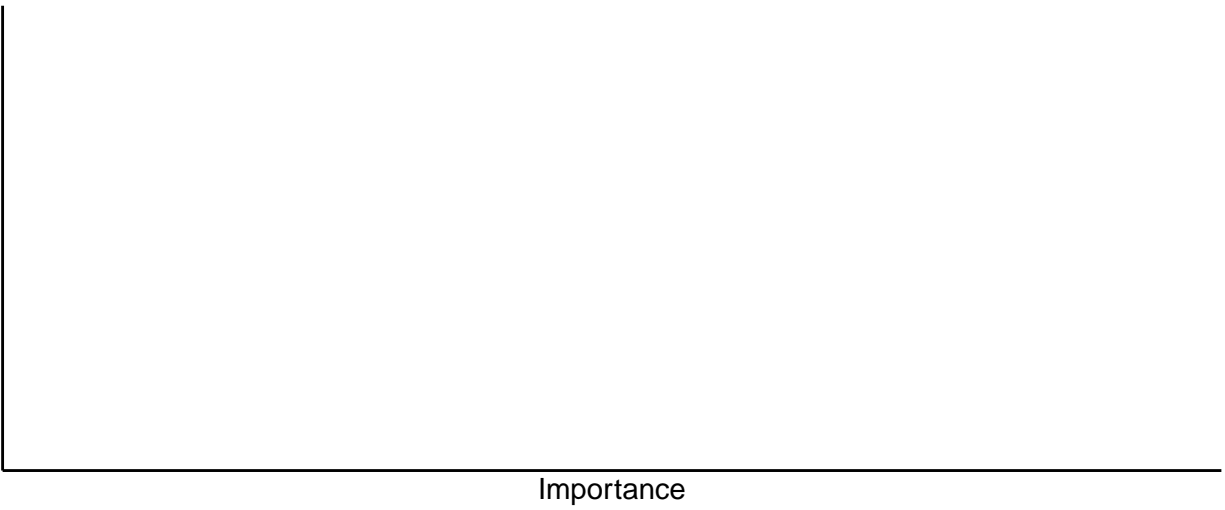
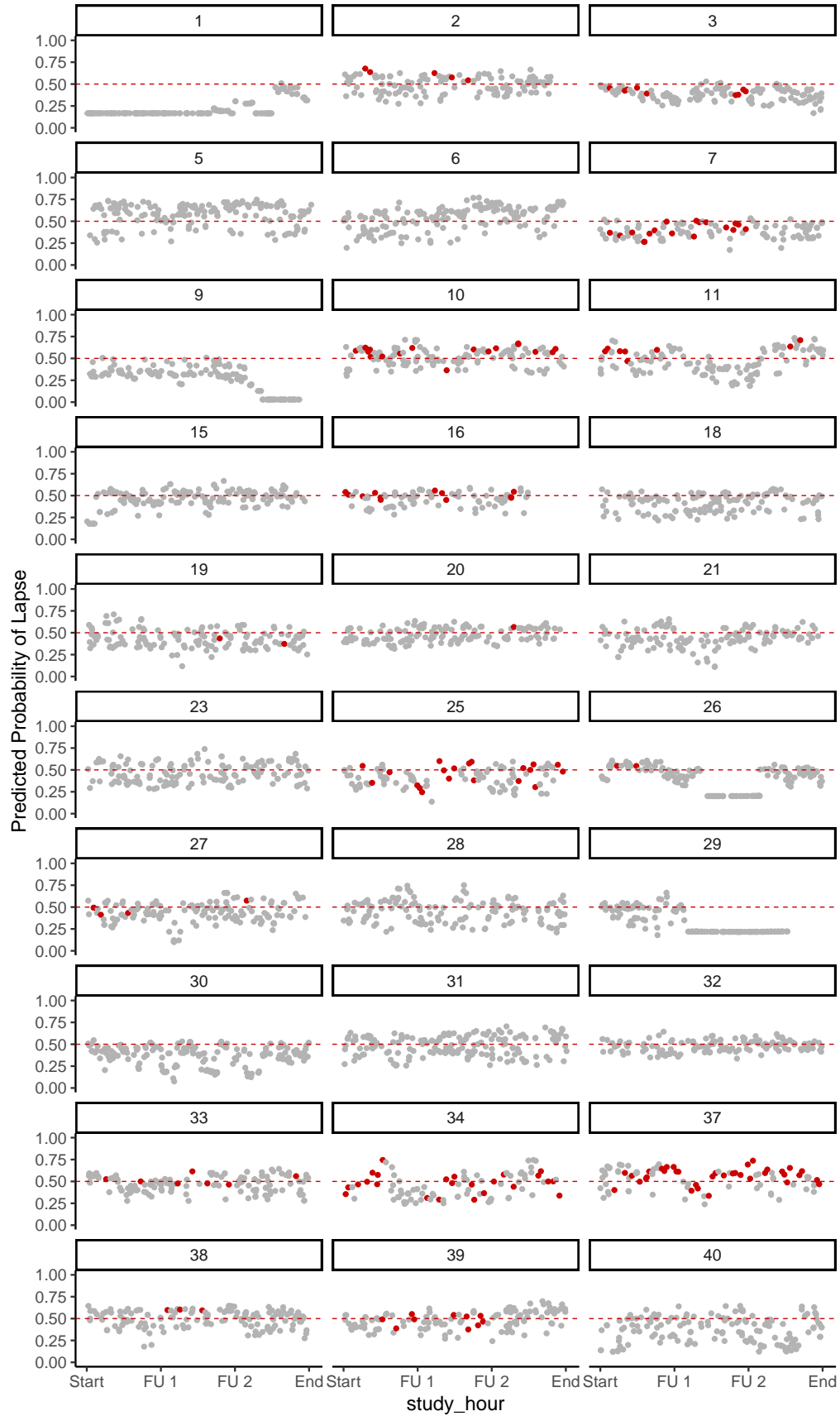
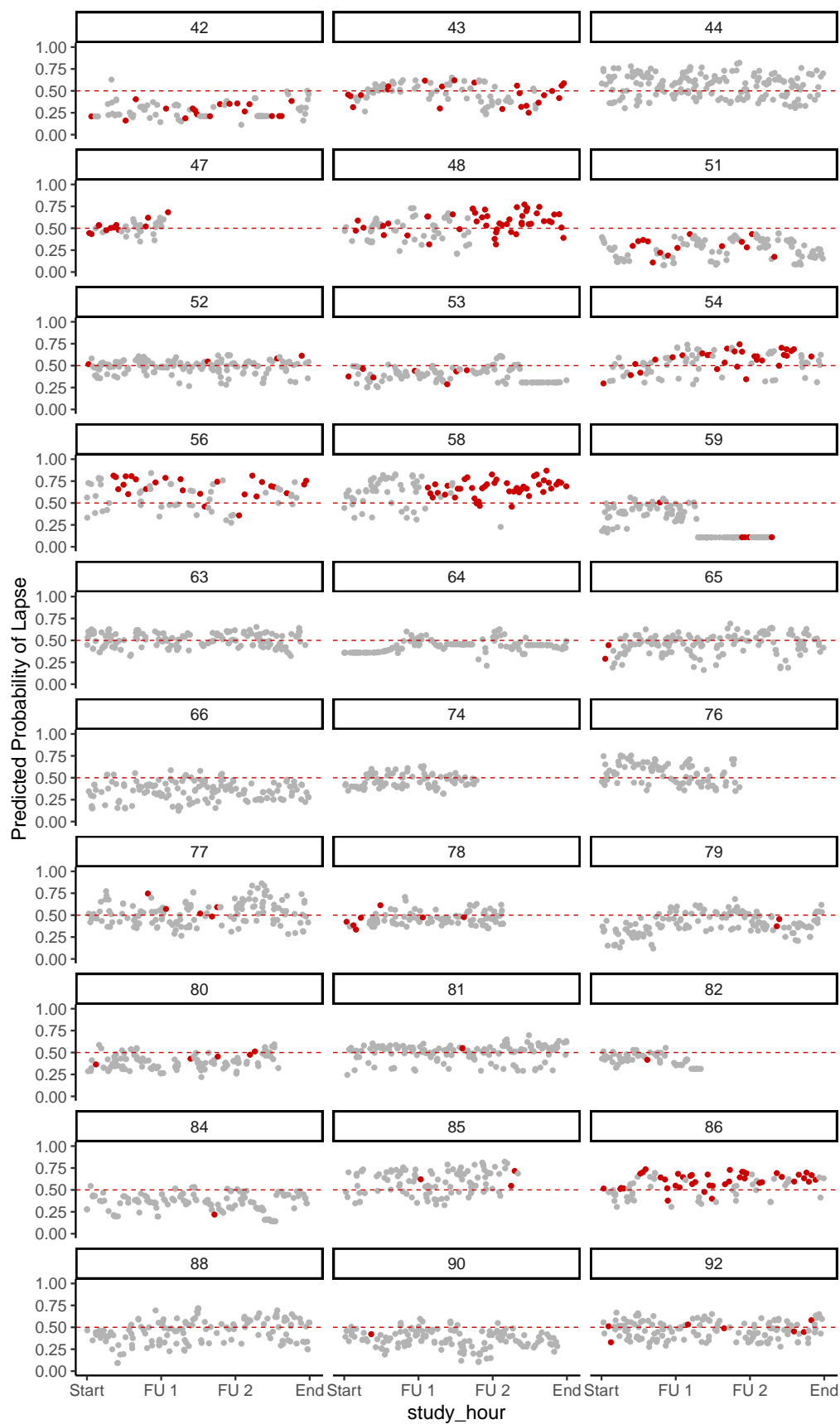
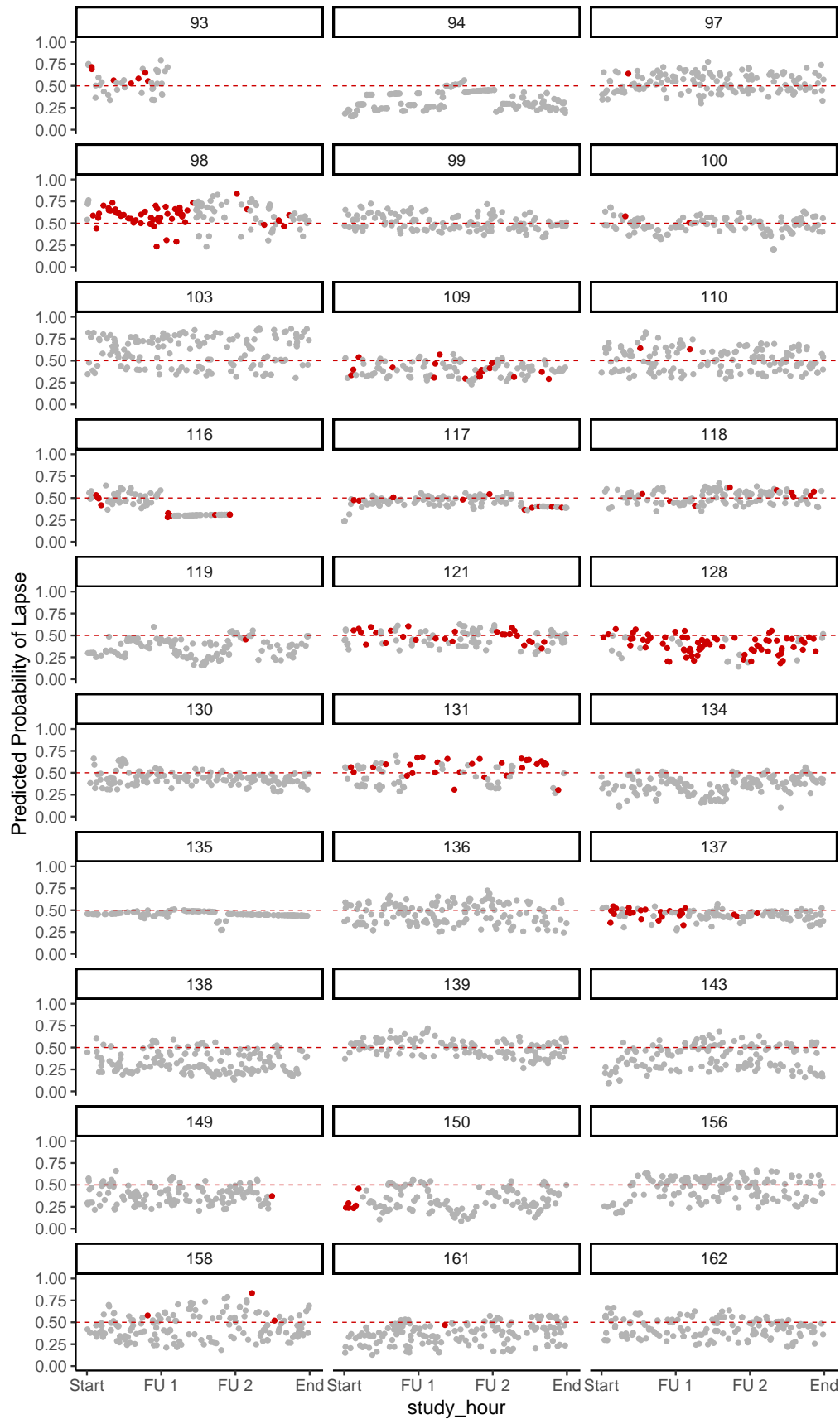


Figure 5. Feature importance scores for best performing active model (glmnet).

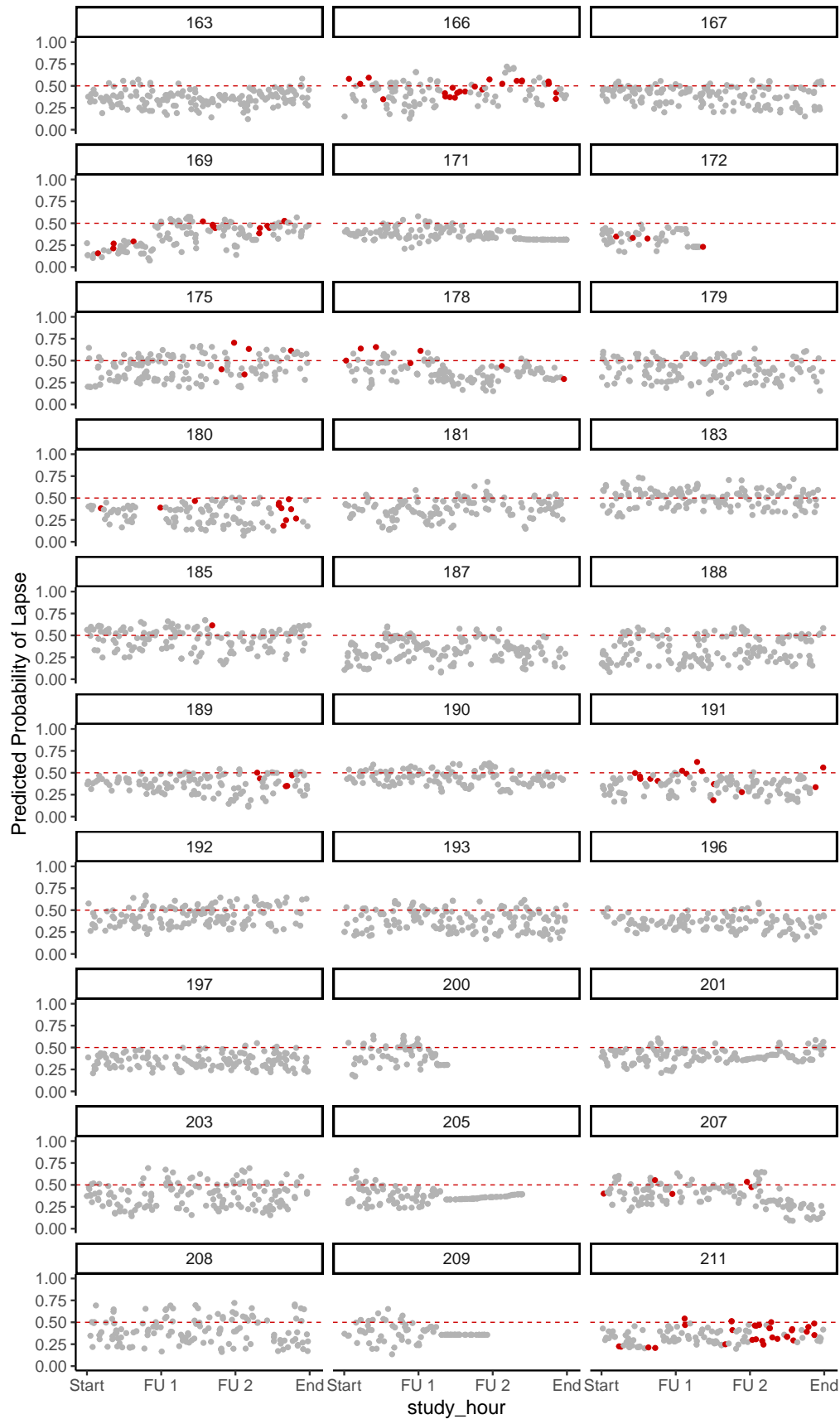
## Appendix

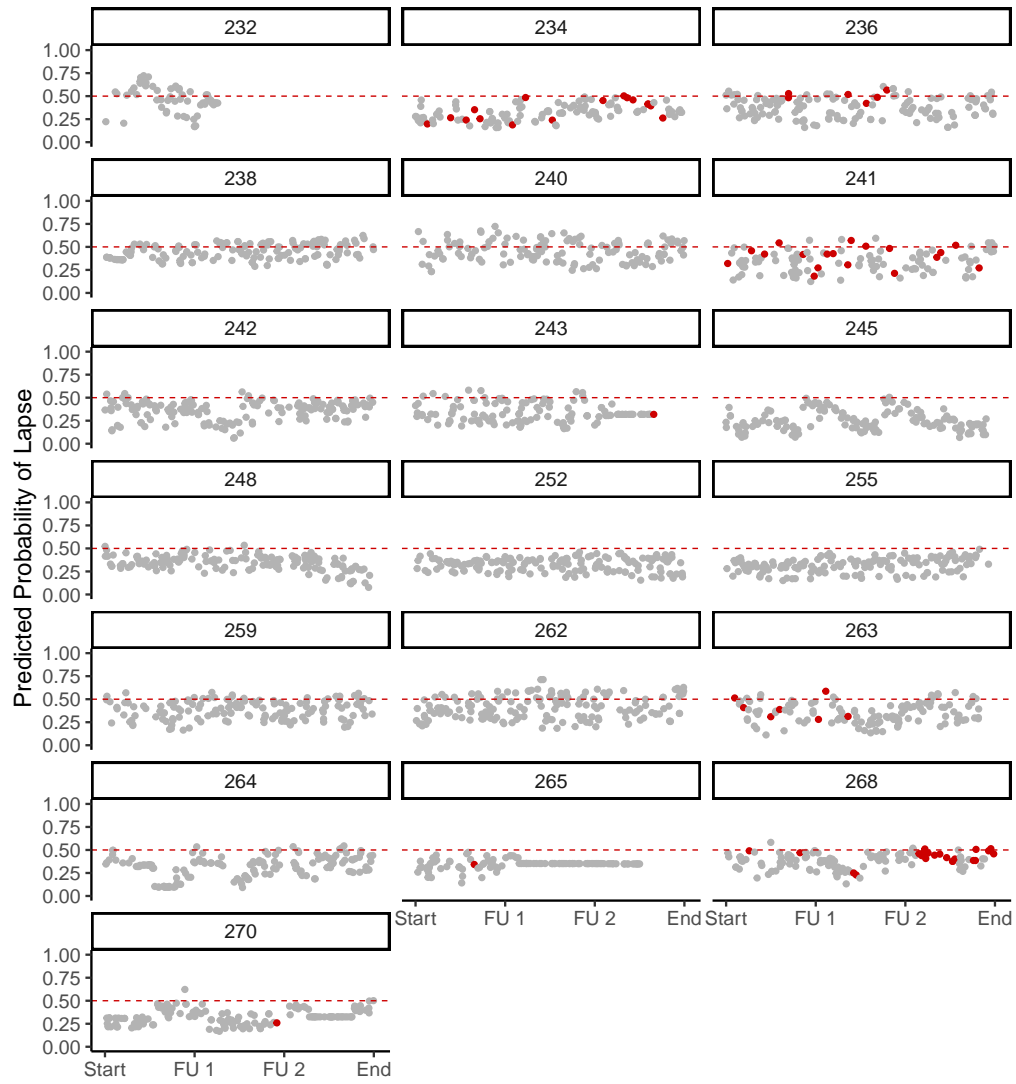












study\_hour

Figure A1. Predicted probabilities of lapse for each participant. A grouped 10x10 resampling method was used to obtain these probabilities. Known lapses are in red. The red dashed line represents the threshold for classifying a probability as a lapse (i.e., everything above the line was predicted to be a lapse).