

Results

Participant Characteristics

Model Selection

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 was a random forest statistical algorithm using passive features. Table 2 characterizes this model over several metrics appropriate for classification. Table 3 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). 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.

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 with no signal. 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.

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 1

```
## # A tibble: 3 x 4
## # Groups:   algorithm [3]
##   algorithm      feature_set      feature_fun_type bal_accuracy
##   <chr>          <chr>          <chr>              <dbl>
## 1 glmnet         feat_all         raw                0.597
## 2 knn            feat_all_passive raw                0.582
## 3 random_forest feat_all_passive perc_raw          0.609
```

Table 2

```
## # A tibble: 7 x 2
##   metric      estimate
##   <chr>        <dbl>
## 1 bal_accuracy 0.598
## 2 sens         0.464
## 3 spec         0.732
## 4 roc_auc       0.643
## 5 accuracy     0.719
## 6 ppv          0.0857
## 7 npv          0.964
```

Table 3

```
##           Truth
## Prediction    no    yes
##           no 143196  5386
##           yes  52314  4904
```

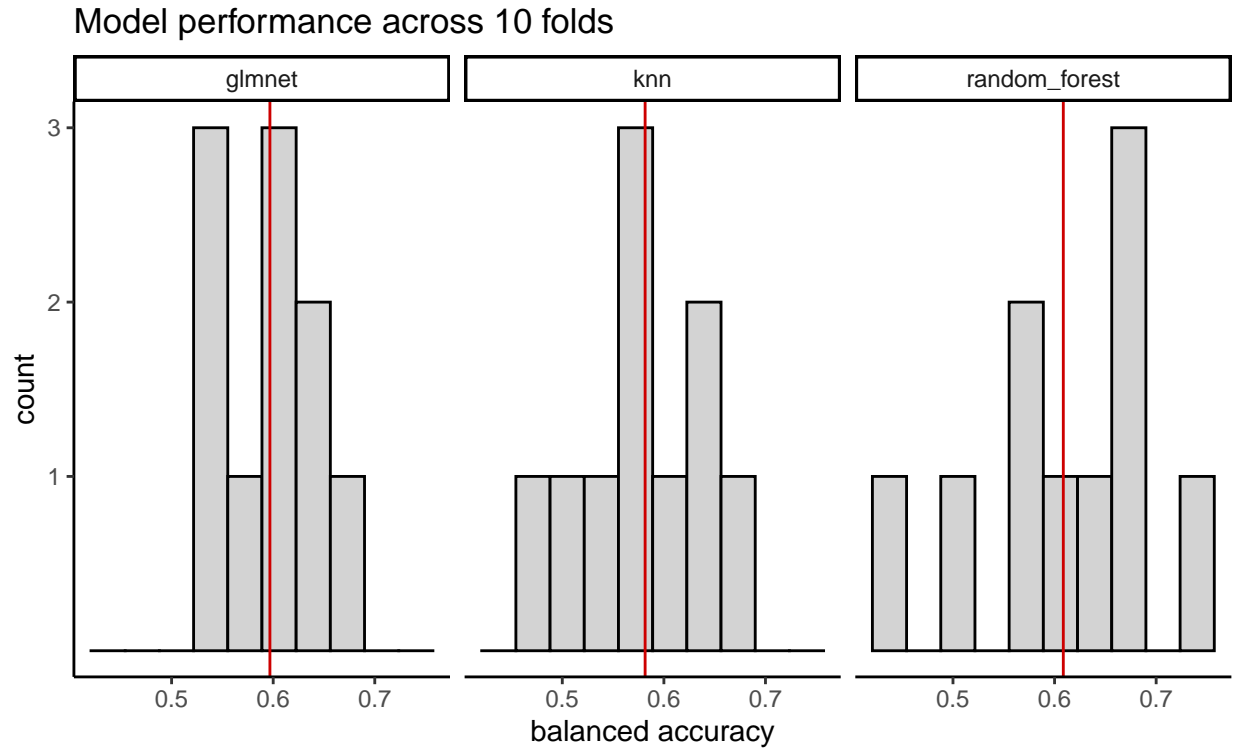


Figure 1. Model performance in each held out fold during model selection.

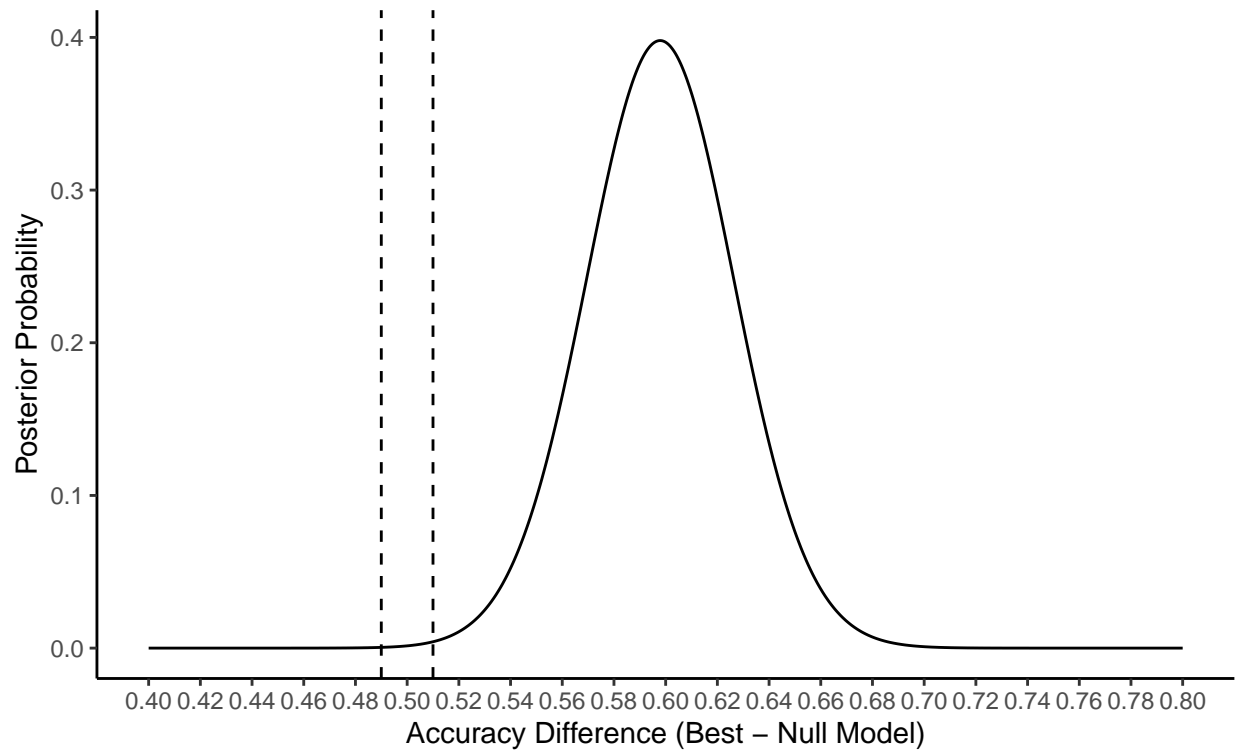


Figure 2. Model comparison of best model and null model.

Context Comparison

Figure 3. Model comparison of passive and active model.

Top Predictors

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.

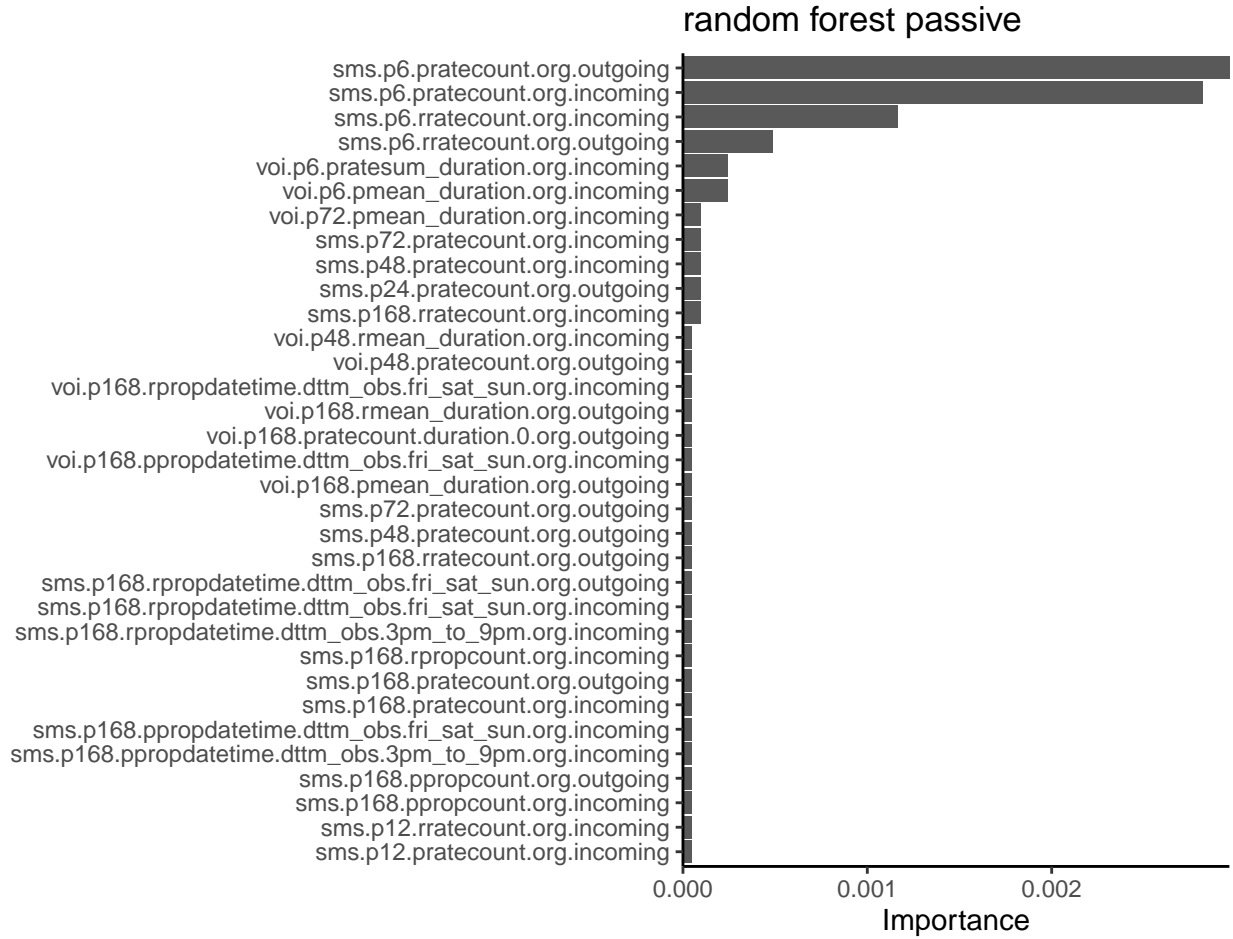
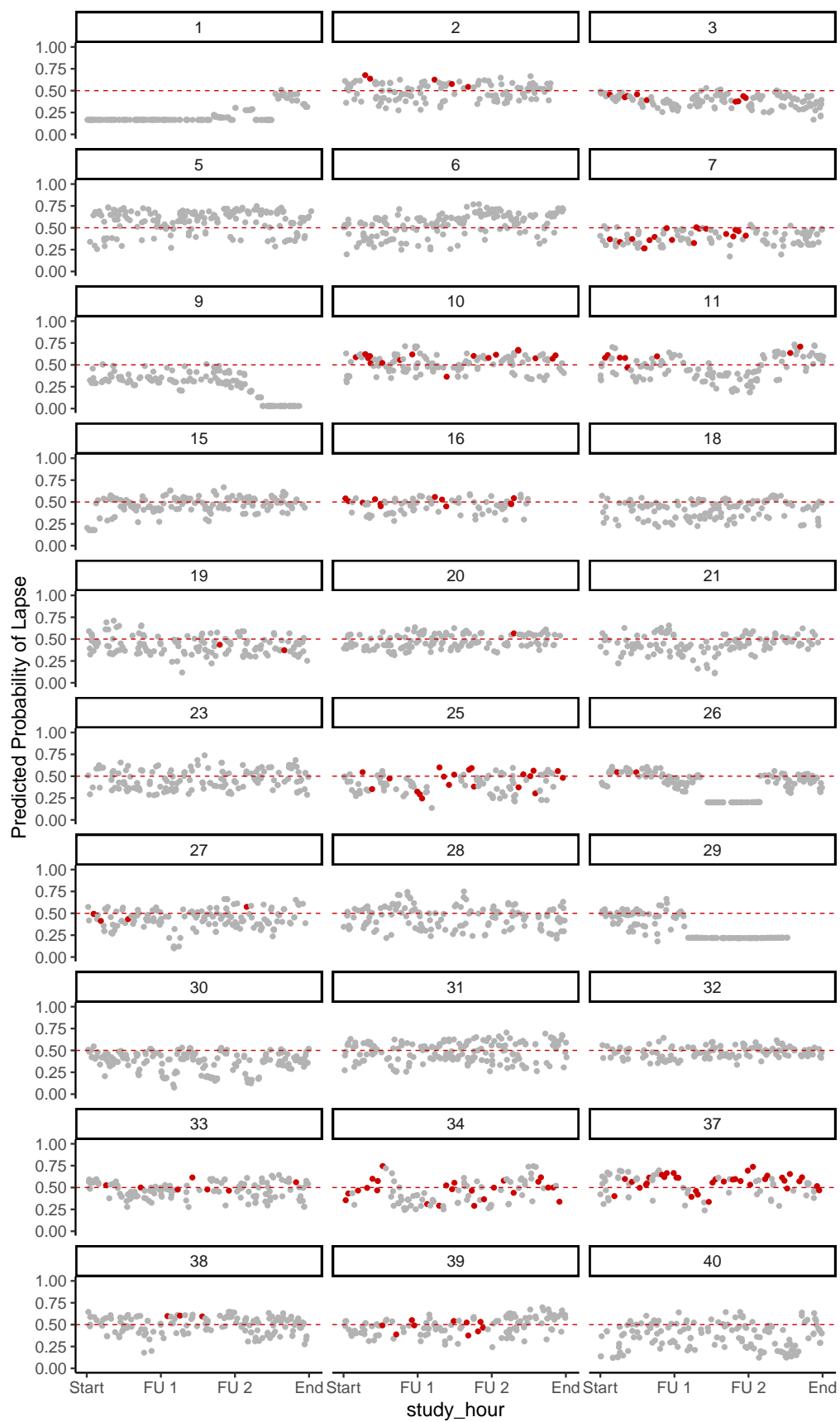
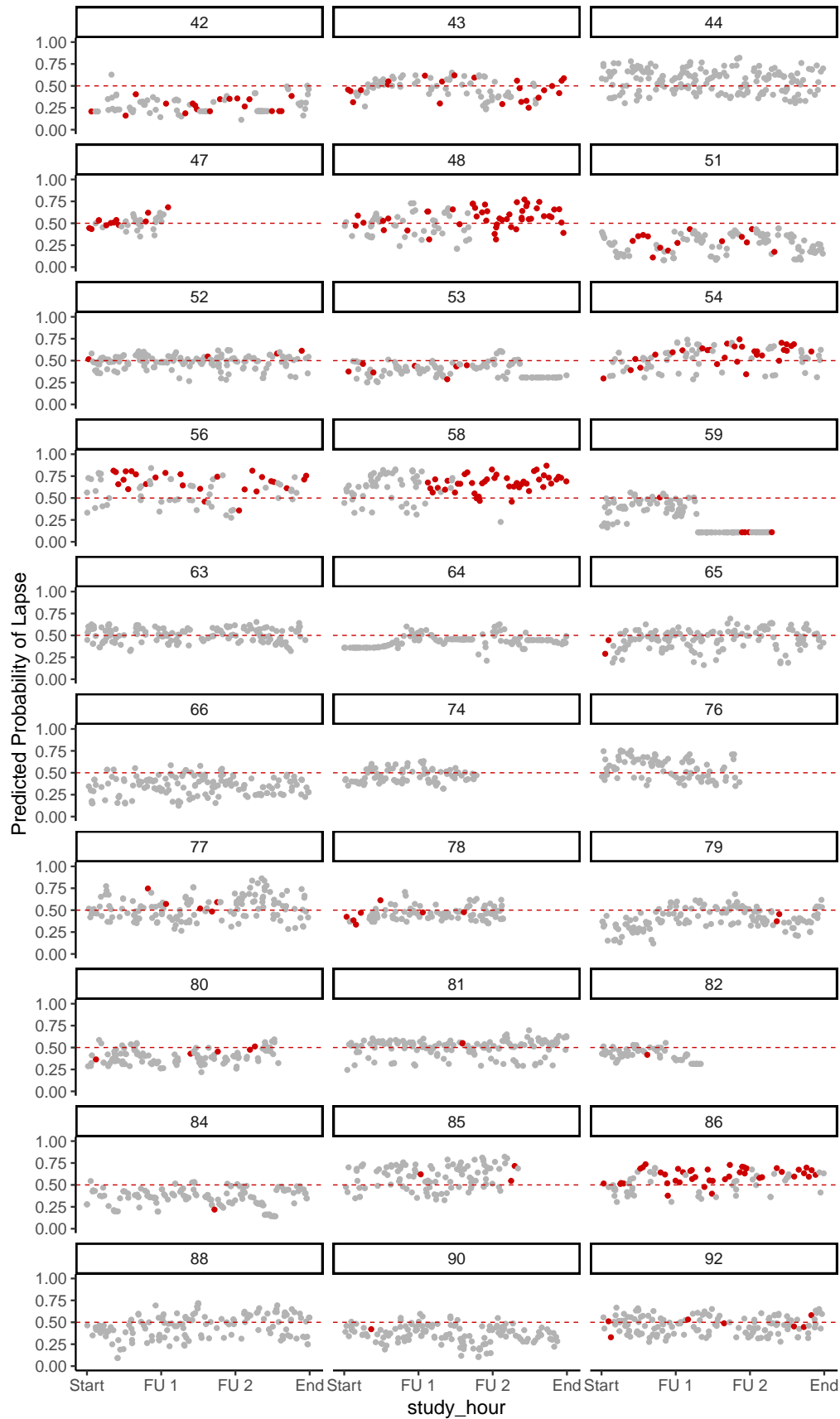
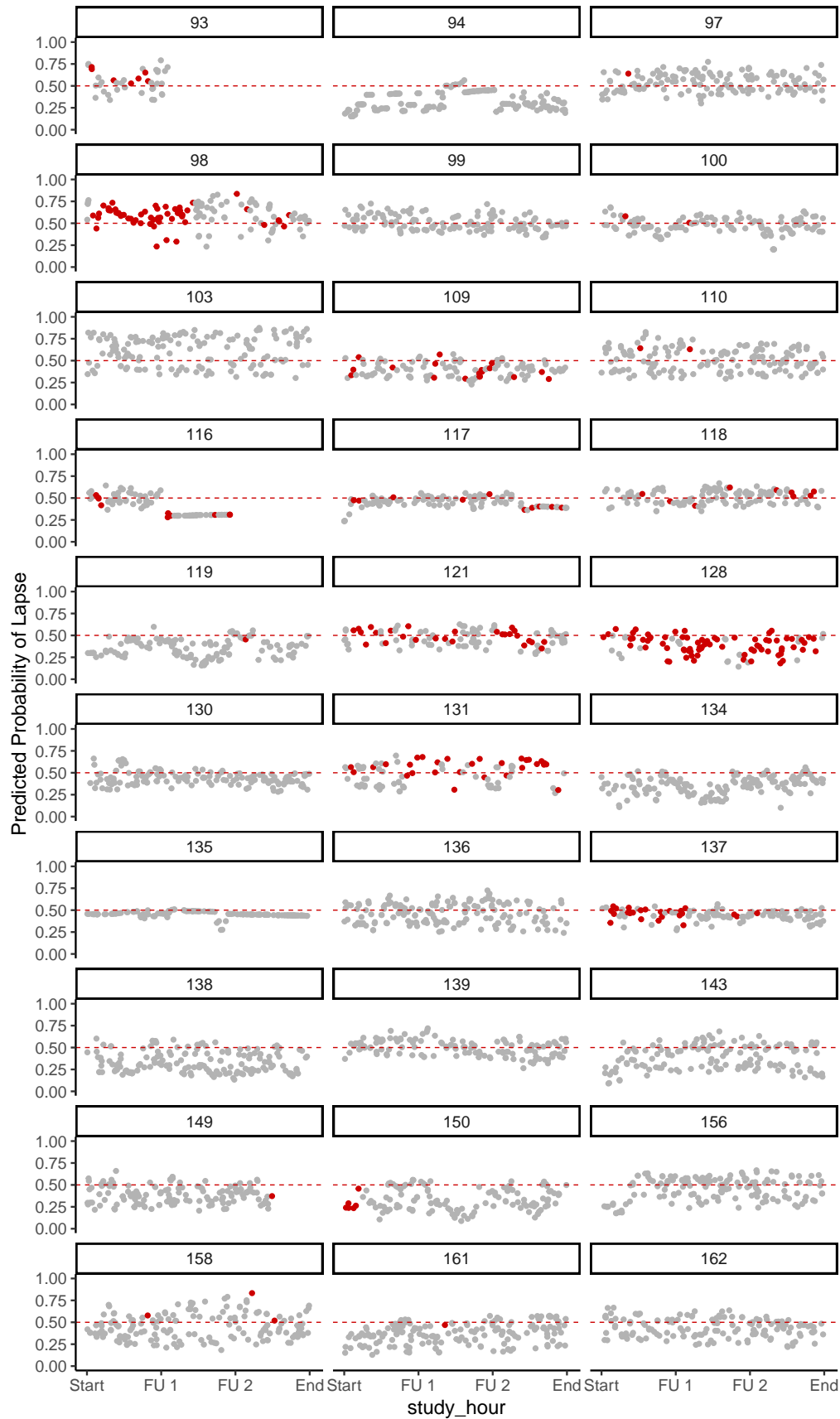


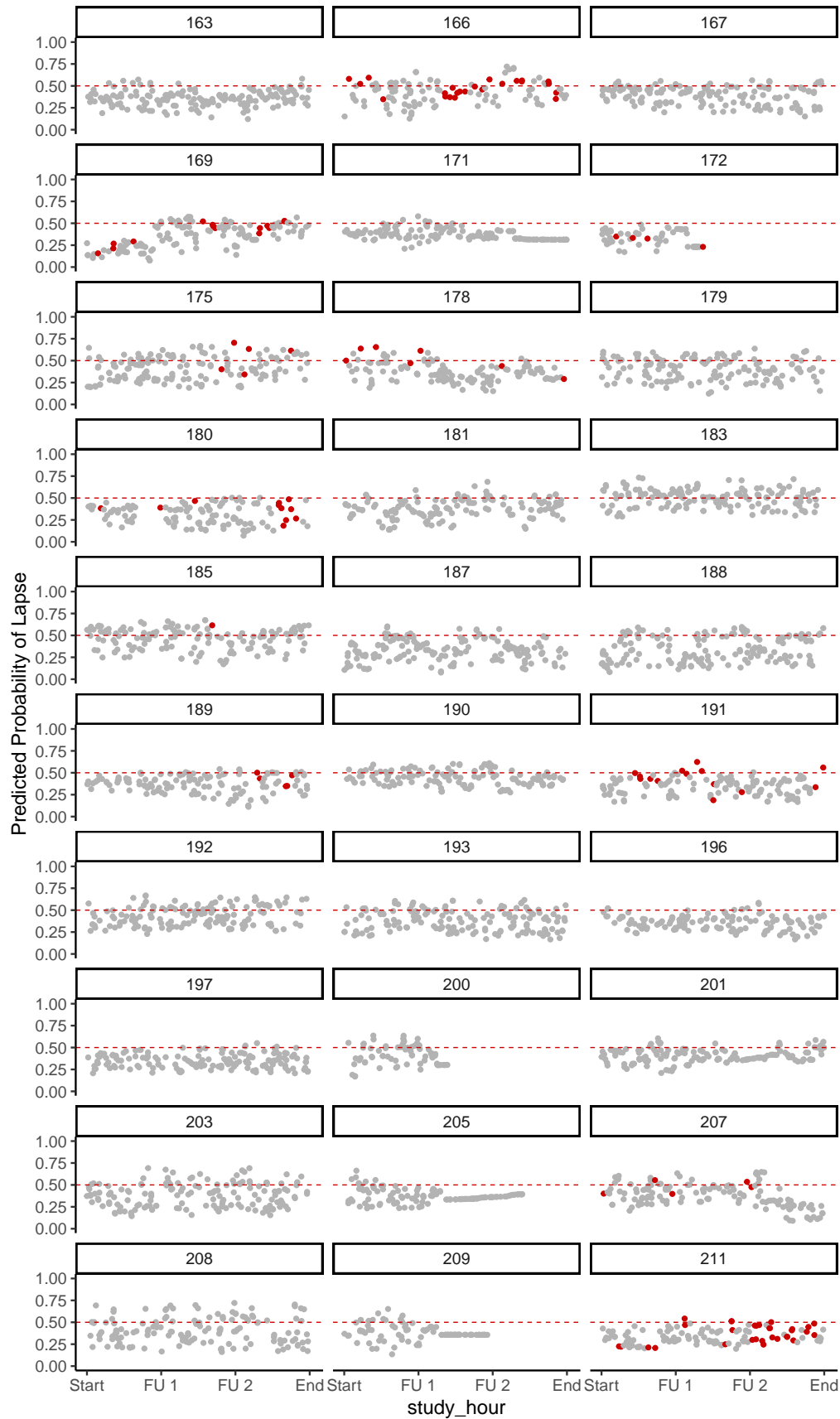
Figure 4. Feature importance scores for best performing model.

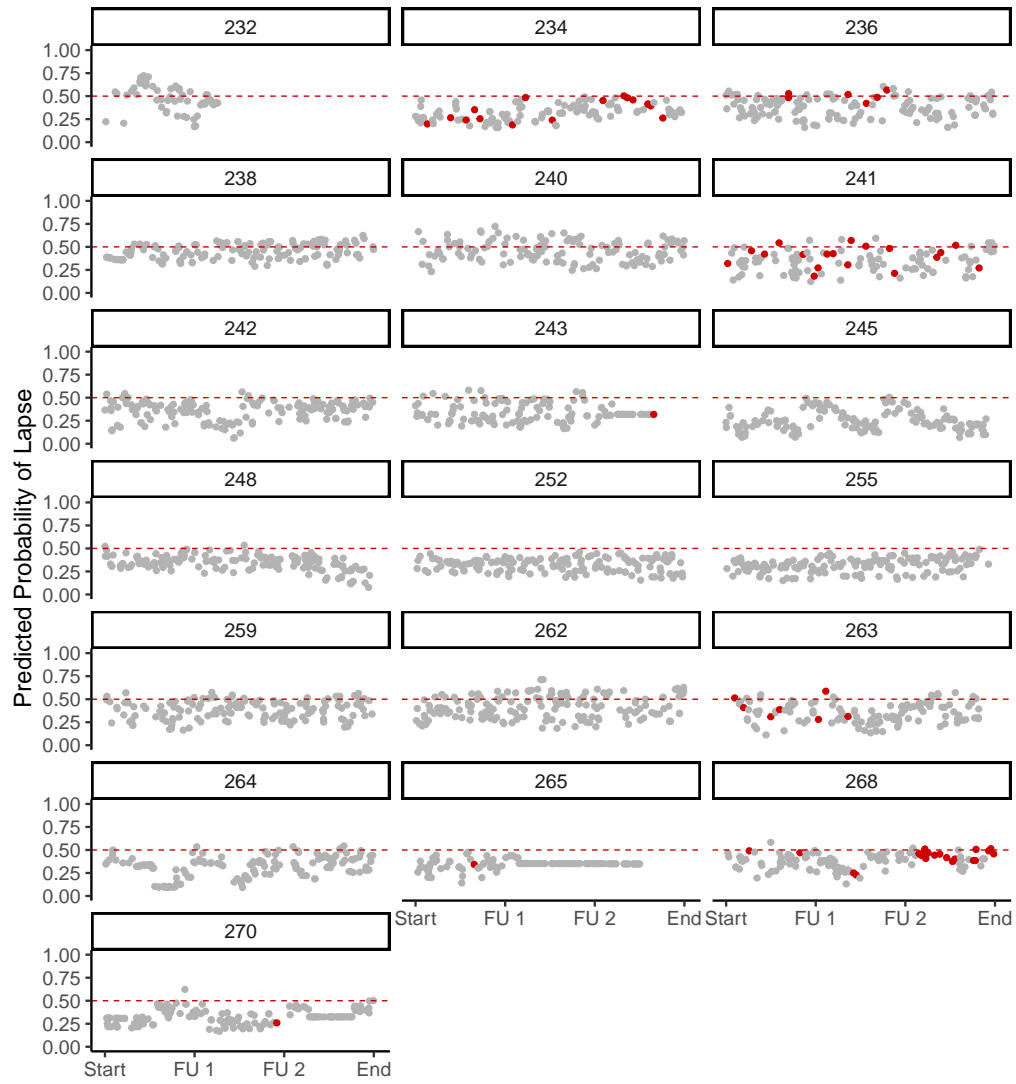
Appendix











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).