Fairbairn, Kang, & Bosch

Supplemental Material

**This material supplements, but does not replace, the peer-reviewed paper in *Drug and Alcohol Dependence***

*Data Processing:* The output provided by the Skyn prototype employed in this research represents the measurement of raw current detected at the sensor and involves no meaningful zero metric. Thus, to approximate a more standardized metric, we subtracted the average of the first five minutes of TAC readings of each session from the entire session as a simple baseline, and one which could be easily implemented in a practical application. [Note that the most recent Skyn prototype also provides measurements in terms of units of alcohol per volume of air, thus providing a standardized metric with a meaningful zero value.] Regarding breathalyzer readings, we obtained regular BrAC measurements from participants in the alcohol condition, but more sparse readings within the no-alcohol condition (see methods section). Participants in both no-alcohol and alcohol conditions were monitored continuously throughout their experimental sessions, and they were further not allowed to keep any possessions with them during their study participation. It was thus possible to infer 0.00% BrAC at times when no alcoholic beverage had been administered by experimenters. Thus, to create instances for the no-alcohol condition, we inserted synthetic (artificial) 0.00% BrAC readings every 10 minutes, so that predictions for no-alcohol participants could be made as well. For the experimental condition, we also added a single synthetic 0.00% baseline reading 1 minute before drinking began in each session. In total, these procedures created 571 instances for the no-alcohol condition and 521 for the alcohol condition.

*Feature Importance Analysis*: We calculated Shapley feature importance values for each instance in the testing set (Lundberg & Lee, 2017). Shapley values describe what a model has learned about the relationships between features and the response variable in terms of the effect each feature had on the prediction for each instance. For example, the value of a feature such as standard deviation of Skyn may have a positive effect on the predicted value in some cases, a negative effect in some, and may have no effect at all in others. Shapley values can be calculated for every instance by tracing the path taken in every decision tree learned by Extra-Trees and recording the influence of the feature on the final prediction. We calculated feature importance by finding the mean absolute Shapley value across all instances for each feature, thus quantifying its total (positive and negative) influence.

We examined mean absolute Shapley values across all instances to discover what features the model found to be the most important in determining the final model predictions. Many features appeared at least once in the model during cross-validation (338), so we examined only the 25 most important (see Figure S1 and Table 1 for expanded descriptions). The three most important were related to the uniqueness of values in the timeseries, which is likely a good indicator of drinking vs. not drinking behavior (if most Skyn readings are identical, it indicates a flat line). Conversely, most of the important features captured change over time in various ways. For example, the four "fast Fourier transform" (FFT) features represent frequency characteristics of the Skyn signal (e.g., repeating patterns), the eight "change quantile" features measure variation restricted to specific quantiles of the data, and the two "linear trend" features capture linear change. Thus, the model appeared to distinguish drinking episodes from non-drinking activity by measuring flat-line Skyn readings, then estimated BrAC during drinking episodes from slope, variation, and frequency-related features.
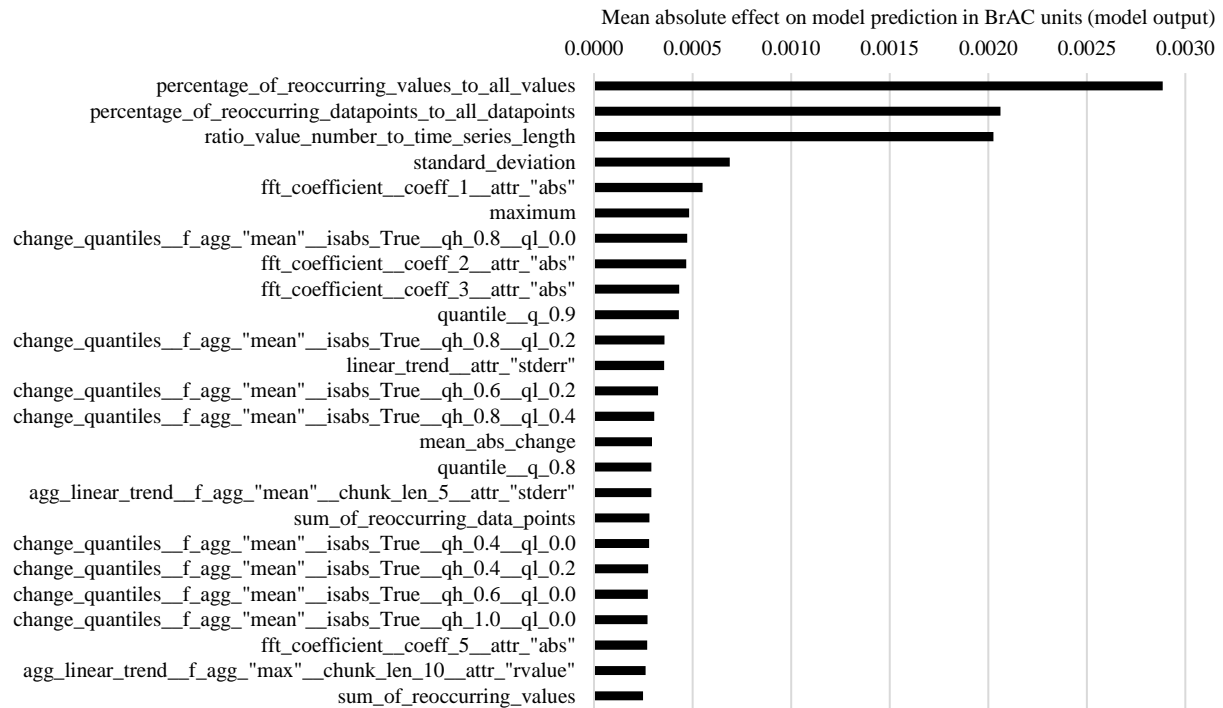
Mean absolute effect on model prediction in BrAC units (model output)

*Figure S1*. Shapley feature importance (mean absolute effect on predicted value) for the top 25 most important features from the main model (TSFRESH features with Extra-Trees machine learning regression). Feature names are from TSFRESH to enable exact matching in TSFRESH documentation; more intuitive descriptions of features are provided in Table 1. In cases where variables may be highly collinear, Extra-Trees will select one variable (essentially at random) when creating each branch in each tree in the model. Thus, total feature importance may be distributed across correlated features.