# BMJ Open

## Digital biomarkers and their potential to increase treatment efficacy in behavioral interventions

SCHOLARONE™
Manuscripts

# Digital biomarkers and their potential to increase treatment efficacy in behavioral interventions

Nicole L Guthrie, MS
Better Therapeutics LLC
San Francisco, California, United States

Jason Carpenter, MS
Manifold, Inc.
Oakland, CA 94607

Katherine L Edwards, FNP-C
Better Therapeutics LLC
San Francisco, California

Kevin J Appelbaum, BSE
Better Therapeutics LLC
San Francisco, California

Sourav Dey, PhD
Manifold, Inc.
Oakland, CA

David M. Eisenberg, MD
Department of Nutrition, Harvard T.H. Chan School of Public Health
Boston, Massachusetts, United States

David L. Katz, MD, MPH
Yale University Prevention Research Center, Griffin Hospital
Derby, CT, United States

Mark A. Berman, MD [corresponding author]
Better Therapeutics LLC
445 Bush Street, Suite 300
San Francisco, California 94108
617-877-0327, mark@bettertherapeutics.io

Word count: 5076

Keywords: digital therapeutics; digital medicine; machine learning, behavioral therapy, hypertension; mobile health; lifestyle medicine

**ABSTRACT**

**Objectives:** Development of digital biomarkers to predict treatment response to a digital behavioral intervention.

**Design:** Machine learning using random forest classifiers on data generated through the use of a digital therapeutic which delivers behavioral therapy to treat cardiometabolic disease. Data from 13 explanatory variables (biometric and engagement in nature) generated in the first 28 days of a 12-week intervention were used to train models. Two levels of response to treatment were predicted: 1) systolic change $\geq$ 10 mmHg (SC model), and 2) shift down to a blood pressure category of elevated or better (ER model). Models were validated using leave-one-out cross-validation and evaluated using area under the curve receiver operating characteristics (AUROC) and specificity-sensitivity. Ability to predict treatment response with a subset of 9 variables, including app use and baseline blood pressure was also tested (models SC-APP and ER-APP).

**Setting:** Data generated through ad libitum use of a digital therapeutic in the United States.

**Participants:** De-identified data from 135 adults with a starting blood pressure $\geq$ 130 / 80, who tracked blood pressure for at least 7 weeks using the digital therapeutic.

**Results:** The SC model had an AUROC of .82 and a sensitivity of 58% at a specificity of 90%. The ER model had an AUROC of .69 and a sensitivity of 32% at a specificity at 91%. Dropping explanatory variables related to blood pressure resulted in an AUROC of .72 with a sensitivity of 42% at a specificity of 90% for the SC-APP model and an AUROC of .53 for the ER-APP model.

**Conclusions:** Machine learning was used to transform data from a digital therapeutic into digital biomarkers that predicted treatment response in individual participants. Digital biomarkers have potential to improve treatment outcomes in a digital behavioral intervention.

**ARTICLE SUMMARY**

**Key messages**

- Digital biomarkers can be created using machine learning to predict therapeutic response and dynamically adjust treatment parameters.
- Predictive models can be used to generate and communicate clinically actionable insights.
- Digital biomarkers have potential to improve treatment outcomes in digitally delivered behavioral interventions.

**Strengths and limitations of this study**

- Practical applications of digital biomarkers in the treatment of cardiometabolic diseases are presented.
- Proof of concept biomarkers demonstrated predictive power despite the small size of the training dataset.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

**INTRODUCTION**

Modifiable behaviors are responsible for 70% or more of all cardiometabolic diseases.[1–3] Health systems are ill equipped to address the current epidemic of cardiometabolic diseases and, in particular, lack widely available behavioral therapies to address the common root causes of these conditions. Digital therapeutics, software designed to encourage changes in behaviors which treat disease, offer a means to deliver behavioral therapy to large populations and preliminary studies demonstrate their potential as a cost-effective treatment for cardiometabolic disease.[4–6]

Compared to pharmacotherapy, digital therapeutics offer potential advantages such as ease of access, ease of use, fewer side-effects and cost-effectiveness.[4, 7–9] Digital therapeutics also generate readily accessible patient data without requiring an office or lab visit. The data generated are voluminous and vary in both type and quality. These can include remotely sensed measures of physiology (e.g., blood pressure, blood glucose, heart rate variability), behavioral data (e.g., about eating, moving, thinking), medication adherence, as well as engagement parameters of unknown significance (e.g., app use, geographic location, circadian patterns of use).

The best use of these data remains an open question. Feeding data directly into EMRs is of limited utility to providers or patients. Whereas, transforming the data into markers of disease status, termed digital biomarkers[10] could provide clinically actionable insights with or without conventional biometric data.[11–13] Digital biomarkers afford a pragmatic approach to remotely monitor patients and intervene on a continuous rather than episodic basis. Greatly expanding opportunities to intervene means that patients have greater access to personalized care, which could improve treatment outcomes.[14, 15]

Machine learning, a type of artificial intelligence (AI) used to make predictions with large and complex datasets, offers a novel approach for creating digital biomarkers. The exponential growth of smartphone use in the United States and advancing interoperability standards allow for digital biomarkers to be compiled across diverse populations and data sources. As a result, the opportunity to advance digital biomarker methodologies has never been greater.[16, 17]

Machine learning is particularly valuable when there is ambiguity about what variables, or to what extent a set of variables predict an outcome of interest. Such ambiguity is inherent to behavioral interventions, like those used in the treatment of cardiometabolic disease. Human behavior results from a complex interaction between the anthropogenic features of our living environment, genetic and epigenetic determinants of behavior, neurobiology, social influences and to some degree chance events.[18–21] While clinical experience and the scientific literature can identify many variables that influence behavior, the interplay of these variables in a given individual and their environment is difficult for a clinician or patient to discern. This ambiguity

limits enthusiasm for behavioral interventions because it makes it difficult for clinicians and patients to rely on behavioral therapy to achieve a predictable level of therapeutic response.

Digital biomarkers can reduce ambiguity by predicting current and forecasting future disease status during the course of treatment.[22–24] Digital biomarkers that serve as markers of current disease status allow for tailoring or adjusting treatment between clinic visits (e.g., when a patient is not doing as well as expected). Similarly, markers of future disease status enable preemptive action, such as adding or subtracting additional treatments, or taking preventive steps to avoid complications of the disease.[13, 24]

In this paper, we present our ongoing work to develop digital biomarkers and aim to illustrate their utility in digitally-delivered behavioral interventions to both the patient and prescribing clinician. We describe the analytic process to show that the development of digital biomarkers requires a hypothesis-driven approach, particularly when datasets are small. Finally, using actual examples of digital biomarkers intended to predict blood pressure status, we discuss the practical and ethical considerations involved in both developing and applying digital biomarkers using machine learning.

**METHODS**
**Digital Therapeutic**
The digital biomarkers discussed in this paper were generated using data from a digital therapeutic created by Better Therapeutics LLC (San Francisco, CA), a developer of prescription digital therapeutics for the treatment of cardiometabolic diseases.

The digital therapeutic integrates a mobile medical application ("app") that delivers behavioral therapy with the support of a remote multidisciplinary care team. The app delivers a personalized behavior change intervention, including tools for goal setting, skill building, self-monitoring, biometric tracking and behavioral feedback designed to provide cognitive training and support the participant's daily efforts to improve overall cardiometabolic disease status. The app facilitates the adoption of evidenced-based behavioral strategies, such as planning and self-monitoring, to increase physical activity and consumption of vegetables, whole grains, fruits, nuts, seeds, beans and other legumes.[25, 26]

Participants were over 18 years of age who self-identified with a diagnosis of hypertension, type 2 diabetes and/or hyperlipidemia. Those reporting to have hypertension at enrollment were offered a free Omron 7 Series Upper Arm Blood Pressure Monitor (Omron Healthcare Inc, Kyoto, Japan) to use during the intervention and to keep at its conclusion. This 12-week intervention was available to all participants at no cost.

This analysis of existing data from the digital therapeutic was approved and overseen by Quorum Review Institutional Review Board[27] and a waiver of informed consent was granted for this retrospective analysis.

**Constraining the Training Dataset**

The clinical intent of this biomarker exploration was to generate digital biomarkers that could improve treatment outcomes in future participants with hypertension for whom blood pressure was not optimally controlled despite the use of pharmacological treatment. We first identified prior participants that met the following criteria: baseline blood pressure was 1) at or above the cutoff for stage I hypertension (systolic ≥ 130 or diastolic ≥ 80)[25] and, 2) recorded no more than 2 weeks prior to or 2 weeks after the start of the intervention.

The development of a predictive model using a training dataset requires all participants to have known outcomes. This is referred to as the "ground truth" in machine learning. In this case, the ground truth was change in blood pressure from baseline. Of the participants identified, we only included those with a follow-up blood pressure value in weeks 7 to 14 of the intervention.

The training set also needs to be a valid representation of future experience to make sure the data sources present are sufficiently representative of those that will be collected in future participants. This is an important consideration because of the evolving nature of digital therapeutics. As the therapeutic is modified over time, the data types collected in earlier versions of the app could be different from those collected in later versions. Therefore, for the purpose of developing biomarkers that predict blood pressure status, we further constrained the training dataset to include participants who used versions of the app which enabled automatic blood pressure tracking.

**Choosing Response Variables and Training Window**

The intent of machine learning is to train an algorithm that can predict a specific outcome, termed the "response variable." The response variable must be both clinically relevant and sufficiently, but not universally, prevalent in the population of interest. For example, if the outcome is "any degree of improvement", but "any degree of improvement" occurs in > 95% of participants in the training dataset, then a predictive model may appear to work well, but could actually be invalid.[28] Furthermore, "any improvement" is arguably less clinically meaningful than "a specific degree of improvement".

To address these concerns, we chose response variables that were well distributed in the training dataset, as detailed in the results section, and were also clinically meaningful in the management of hypertension. Specifically, we chose to predict: 1) a systolic blood pressure improvement of 10 mmHg or higher near the end of a 12-week intervention period (defined as 7 to 14 weeks after start), because 10 mmHg has been well demonstrated to be clinically meaningful[29, 30] and that

degree of change is near the mean for participants in the training dataset (66/135 participants); and 2) a reduction in blood pressure to the elevated range or below (systolic ≤ 130), because this level of blood pressure control would signal a clinician to stop adding additional pharmaceuticals or consider reducing or deprescribing pharmaceuticals (48/135 participants).[25]

For a digital biomarker that predicts blood pressure status to be clinically useful, it needs to compute with data collected in a short period of time, and in less time than typically occurs between clinic visits. This means the biomarker could be used to intervene between office visits and could play a role in addressing clinical inertia that limits primary care providers ability to optimize blood pressure control in their patients.[31, 32] To demonstrate proof of concept, we chose a 28-day training interval, meaning that we trained machine learning models on the first 28 days of patient data to evaluate whether it could predict blood pressure change in weeks 7 to 14. We hypothesized that data collected within this short training window could sufficiently represent changing behavioral patterns and treatment response, so as to predict future blood pressure status.

**Choosing Modeling Techniques and Explanatory Variables**

There are many valid methodologies for machine learning. These methods can be categorized by the type of learning involved - supervised or unsupervised. While a discussion on the merits and nature of each type is beyond the scope of this paper, it is important to select modeling techniques appropriate for the size of the dataset and nature of response variables chosen. For the biomarkers presented herein, we utilized random forest models, which is a form of supervised classification learning known to reduce overfitting in small datasets.[33] The models are supervised because the ground truth (i.e., actual blood pressure change) for each participant in the training dataset is labeled. The models are attempting to learn whether the classification was or was not achieved; for example, will a participant achieve a > 10 mmHg blood pressure reduction, or not?

In addition to classification labels, random forest models must be trained on a set of explanatory variables. For a small training dataset, each model must have a limited number of variables so as to avoid excess noise and overfitting that can lead to reduced generalizability. These explanatory variables must be selected by hypothesis. For example, we hypothesized that baseline blood pressure and achievement of behavioral goals would influence the degree of blood pressure change observed and used these as explanatory variables. In a large dataset, feature engineering can be used to identify the most predictive explanatory variables.

For each biomarker model, denoted SC (systolic change of 10 mmHg or more) or ER (elevated range achieved), we used 13 explanatory variables, which can be categorized as engagement or biometric variables. Engagement variables are counts of actions related to the use of the digital therapeutic, including count of all meals reported, plant-based meals reported, physical activity

reported and length of exposure to the intervention. Biometric variables included baseline systolic, baseline diastolic, mean systolic and diastolic at training window end, initial systolic and diastolic change (end training mean - baseline), minutes of physical activity, and baseline Body Mass Index (BMI).

We trained each biomarker model on data generated during a 28-day training window, starting with participant day 0 (sign up) up to day 28 (a third of the way through the studied 12-week intervention period). We utilized hyperparameter optimization to minimize overfitting and to achieve the maximal leave-one-out cross-validated area under the receiver operating characteristic curve for both models.

**Model Validation**

Performance of each biomarker model was assessed using leave-one-out cross-validation, which is a common and valid technique for use in samples of this size.[34–36] This was done by training each model on N-1 samples of the data and then making a prediction on that one sample that was left out, producing an "out of sample' prediction for all N samples. The N predictions were pooled to generate the classification variables of the receiver operator characteristic curve (ROC), the area under the curve of the ROC (AUROC) and a confusion matrix of true versus predicted values.[37]

For each biomarker model, the ROC curve illustrates predictive ability of the response variable (in this case systolic change of 10 mmHg or move to a range of elevated BP) at different thresholds of discrimination. At each prediction discrimination threshold, the ROC displays the false positive rate (FPR) against the true positive rate (TPR). The FPR is the ratio of truly negative events categorized as positive (FP) to the total number of actual negative events (N). Specificity or true negative rate of a model is calculated as 1 - FPR and is an indication of how well a model does in correctly identifying those who do not achieve a successful outcome, as defined by the response variable.

Since the intended application of these biomarkers is analogous to a diagnostic test, which are traditionally evaluated based on their specificity, we evaluated model performance at a specificity of 90% (FPR = .10). A low FPR minimizes the number of participants who the model would predict to achieve a healthier state who actually will not. In turn, this minimizes the number of participants who might be erroneously taken off blood pressure medicine as a result of an erroneous prediction. It is less critical to avoid labeling participants who had achieved a healthier state as though they had not. This is why we choose a discrimination threshold with a low FPR, and then evaluate the TPR at that point.

As a further validation step, we examined the performance of each biomarker by excluding the four explanatory variables that capture blood pressure change in the training window. While we

hypothesized that these models would perform less well, they serve to test the concept that a digital biomarker that predicts blood pressure status can be generated without using ongoing blood pressure data. These validation models using only app engagement and other biometric variables, are denoted SC-APP and ER-APP.

**Making Use of Explainable Artificial Intelligence (AI)**

Digital biomarkers that are generated using machine learning do not need to be viewed as a black box. Instead, explainable AI techniques are available that can provide more granular details about the explanatory variables that influenced the prediction made. Explainable AI can afford both individual participant and population level insights.

We used the Tree Shapley Additive Explanation (SHAP) algorithm on the random forest models[38] to generate more interpretable predictions at the participant level. The SHAP algorithm assigns each explanatory variable an importance value for each prediction. Using SHAP on a machine learning model is analogous to coefficient analysis in classical regression. Similar to coefficient analysis, it can be used to determine the relative importance of explanatory variables in addition to determining which explanatory variables drove a particular prediction. Predictions start at a base value that is the expectation of the response variable. For binary classification models, this is defined by the proportion of outcomes by class (e.g., the proportion of participants who successfully reduced their blood pressure). Then SHAP values attribute to each explanatory variable the change in expected model prediction given the addition of that explanatory variable. This provides insight into how much each explanatory variable positively or negatively impacts the prediction made for each participant. The final prediction probability of whether the participant will achieve the response variable is the sum of the base value and all of the explanatory variable attributions.

Individual SHAP values can be used to provide specific behavioral feedback to participants, with the intent of motivating a change in behavioral pattern that may improve treatment outcomes. SHAP values for all participants can also be plotted to reveal the overall ranking of variables in the population studied. These variable rank lists can then inform hypotheses about how to further improve the design of the digital therapeutic to optimize clinical outcomes.

All machine learning model development was done using open-source packages in Python. The packages include but are not limited to Scikit-Learn, SHAP, Pandas, and Numpy.

**Patient and Public Involvement**

We did not directly involve patients or public (PPI) in the current study. However, the digital therapeutic that is the source of the database used in this study was developed with PPI input over the product development lifecycle.

**RESULTS**

**Dataset**

The training dataset contained 135 participants who met the inclusion criteria. The mean age was 54.9 years (95% CI 53.5, 56.3), mean baseline BMI was 34.5 (95% CI 33.1, 35.8) and 83% (112/135) were female. Based on the 2017 ACC/AHA definition,[25] half of the participants (68/135) had stage 1 hypertension at baseline, with the other half (67/135) having stage II hypertension at baseline. Of those with stage 1 hypertension at baseline, 51.5% (35/68) had isolated diastolic hypertension (i.e., diastolic BP 80-90 mmHg). Of those with stage II hypertension at baseline, 14.9% (10/67) had isolated diastolic hypertension (i.e., diastolic BP ≥ 89 mmHg). Baseline characteristics are listed in Table 1.

Over the intervention period examined, systolic blood pressure changed by -12.7 mmHg (95% CI -14.8, -9.6), diastolic blood pressure changed by -7.1 mmHg (95% CI -9.0, -5.2) and the mean duration of days between baseline and most recent value for blood pressure was 79.3 days (95% CI 76.8, 81.9). Of all participants, 35.6% (48/135) shifted to a blood pressure range below stage I (<130/80). A shift to a normal range was seen in 16.4% (11/67) of those starting with stage II hypertension and 29.4% (20/68) of those starting with stage I hypertension.

**Table 1**. Participant characteristics at baseline

| Participant Characteristics | Total (n = 135) | Stage I BP (n = 68) | Stage II BP (n = 67) |
|---|---|---|---|
| | | | |
| Age (years), mean (95% CI) | 54.9 (53.5 to 56.3) | 55.7 (53.7 to 57.7) | 54.2 (52.1 to 56.2) |
| Body mass index (kg/m$^2$), mean (95% CI) | 34.5 (33.1 to 35.8) | 33.8 (31.9 to 35.6) | 35.2 (33.2 to 37.2) |
| Female, n (%) | 112 (83) | 54 (79.4) | 58 (86.6) |
| Systolic BP (mmHg), mean (95% CI) | 138.9 (136.2 to 141.6) | 127.9 (126.1 to 129.7) | 150.0 (146.6 to 153.5) |
| Diastolic BP (mmHg), mean (95% CI) | 87.8 (86.1 to 89.4) | 82.3 (81.1 to 83.6) | 93.3 (90.7 to 95.8) |

| Isolated diastolic hypertension, n (%) | 45 (33.3) | 35 (51.5) | 10 (14.9) |
|---|---|---|---|
| BP medications (count), mean (95% CI) | 1.3 (1.1 to 1.5) | 1.2 (1.0 to 1.4) | 1.4 (1.1 to 1.7) |

**Predictive Models**

The random forest classifier achieved optimal performance with 100 trees and a minimum of 3 samples per leaf node for the SC model. For the ER model optimal performance was achieved with 400 trees and a minimum of 5 samples per leaf nodes.

Biomarker models were assessed at the operating point on each ROC that was as close as possible to a FPR of 10%. The SC model (predicting a systolic change of 10 points) was assessed at a FPR of 10%, which means that 10% of participants who didn't achieve a reduction in systolic blood pressure of 10 mmHg were labeled as though they had. Evaluating the model at 10% FPR, we were able to achieve a TPR of 58%. This means that 58% of participants who achieved a reduction in systolic blood pressure of 10 mmHg were labeled correctly. The AUROC was .82, model specificity (1 - FPR) was 90%, sensitivity (TPR) was 58% and accuracy ((TP+TN)/n) was 74%. In the SC-APP model, where variables related to changes in blood pressure were removed, the AUROC was .72 and at a FPR of 10% (specificity of 90%), the TPR was 42%. The resultant receiver operator curves for these 2 models can be seen in Figure 1.

The biomarker models exploring the ability to predict a shift down to a blood pressure range of elevated or better (ER and ER-APP) also demonstrated predictive capacity, but less so than the SC models. For the ER model, the AUROC was .69 and at a FPR of 9% the TPR was 32%. When the BP change variables were removed, in the ER-APP model, the prediction ability was only slightly above chance (AUC = .53, TPR 26% at FPR of 12%).

Plots of the Tree SHAP algorithm results for the SC and SC-APP models are shown in Figure 2. Explanatory variables on the y-axis are ordered from most to least predictive based on their average absolute contribution to the response variable. Each dot represents the SHAP value of that variable for one participant and the placement of the dots on the x-axis indicate if the contribution was subtractive or additive for a specific participant. The color of the dot is indicative of the value for that variable, with highly positive values displayed as red and low or negative values showing up as light blue. The plot for the SC model reveals that explanatory variables related to blood pressure were top contributors to the prediction. For example, the distribution of dots across the x-axis for the first variable listed shows that improvements in systolic BP early in the intervention, as seen by the blue and purple dots to the right of 0 on the

x-axis, contributed positively to the prediction that a participant would succeed. Behavioral variables also had predictive power. For example, a high count of physical activity minutes and plant-based meals reported positively contributed to a prediction of success for most participants.

Shapley values can be aggregated and illustrated for every participant. A plot of the SHAP values helps to visualize which variables contributed most to a low or high prediction of success for an individual participant. In figure 3, we display the SHAP values for two participants, one with a lower than expected probability of success (example A), and one with a higher than expected probability of success (example B). In example A, the participant experienced a large improvement in their systolic blood pressure in weeks 3 and 4 (-14 mmHg), yet is given a low probability of sustaining this improvement at the end of the intervention period. This surprisingly low probability is explained by the SHAP values, which reveal low counts for several behavioral explanatory variables, such as the number of plant-based meals and minutes of physical activity reported. This data can be automatically translated into a simple explanation to the participant, that their probability for sustaining meaningful change could be higher if they made incremental improvements in their meal and activity pattern. Furthermore, the exact number of additional meals and activity minutes to accrue per week to sufficiently increase the probability of success could be calculated, to motivate the participant and to give meaning to the additional efforts prescribed.

In example B, the participant has evidenced no improvement in systolic blood pressure at the end of week 4, yet they are predicted to meaningfully improve blood pressure by the of the treatment period. This unexpected prediction is explained by the SHAP values, which show that the combined impact of their baseline blood pressure and behavioral explanatory variables suggests a high likelihood of success. This data can be automatically translated to provide timely encouragement to the participant to maintain or advance their behavioral changes even though their blood pressure has not yet responded.

**DISCUSSION**

In this paper, we present a proof of concept that digital biomarkers can be developed using even a small training dataset from a digital therapeutic. We demonstrated that 28 days of data can be transformed, using machine learning, into a digital biomarker that predicts the degree of treatment response, in this case whether a meaningful drop in blood pressure will occur at the end of the treatment period.

There are many ways to use such a biomarker in practice to tailor behavioral treatment and improve outcomes. For a patient, these biomarkers can be used as a continuous form of treatment feedback and behavioral reinforcement. The probability of a significant treatment response can be translated into a treatment score, much like a credit score. Since this score could be recalculated with every new engagement recorded in the digital therapeutic, it would serve to

motivate app use and reinforce healing behaviors. In addition, the biomarker output can be made more meaningful using explainable AI techniques. For example, SHAP values can be translated into a prioritized list of behavioral actions to help a patient focus their attention on efforts that are most predictive of success.

For clinicians and health systems, digital biomarkers can function as a form of automated patient monitoring. The probability of a positive treatment response can be translated into a clinical alert by setting an acceptable specificity-sensitivity threshold for each biomarker paired with a duration of time above this threshold. Like any diagnostic test, the performance characteristics of the alert should be made known to those acting upon it. Since such an alert would be intended to influence treatment decisions, for example via a clinical decision support tool, specificity-sensitivity pairs need to be evaluated from a risk-benefit perspective. For instance, how do the risks associated with false positives and false negatives compare to the benefits of identifying true positives and true negatives? To accurately weigh these risks and benefits requires us to understand the context that the biomarker and therapeutics are used. Current clinical practice, for example, is plagued by high rates of clinical inertia (i.e., a lack of timely and appropriate treatment decisions). Therefore, a higher false positive rate may be tolerated as a trade-off for an easier-to-access biomarker.

For a developer of digital therapeutics, digital biomarkers provide not just a way to personalize treatment and communicate clinical status to providers, but also a way to better understand what variables within the therapeutic are most predictive of clinical outcomes. These data can be used to guide the ongoing refinement of a digital therapeutic. When datasets are of sufficient size, the machine learning techniques used to generate digital biomarkers can also be applied to identify distinct digital phenotypes, that is, unique patterns of engagement with a behavioral intervention that represent meaningful subpopulations who share the same diagnosis. Identifying and targeting treatment to previously unknown subpopulations is thought of as meaningful step towards more personalized medicine.[15, 39]

**Limitations, Practical and Ethical Considerations**
The main limitation of the work presented here is the size of the training dataset used. It is likely that a larger dataset would improve the performance characteristics of biomarkers tested.[14] It is noteworthy, given this limitation, that one of the biomarkers (SC) had an AUC greater than 0.8. This suggests the utility of beginning digital biomarker development in early in the implementation of a digital therapeutic.

Other limitations include the omission of known predictors of treatment response (e.g., medication adherence or change), the reliance on a small set of explanatory variables and the inclusion of self-reported variables that may be subject to human error. Addressing these limitations may enable more accurate biomarkers.

The ethical and practical implications of applying complex, ever-changing, predictive models generated by machine learning are only beginning to be appreciated. To preempt potential misuse of digital biomarkers, we must understand the true meaning of the data these biomarkers present. For example, a predictive model can identify variables that are predictive of a given outcome. This does not mean highly predictive variables caused the outcome, nor does it mean that poorly predictive variables are not causative. Instead, the predictive strength of each variable should be treated literally as "markers." This does not preclude the automated use of explanatory variables to guide the personalization of behavioral therapy. However, since the individual level of each variable could be influenced by unknown confounding factors, and since the degree of modifiability of each variable is not known, the impact of this form of behavioral therapy must be studied.

Finally, like any biomarker, a digital biomarker is only generalizable if the training dataset is truly representative of future patients. If the training dataset is biased or overly skewed, it may produce a biomarker that underperforms at best, and is harmful at worst. To guard against this bias, re-validation of a digital biomarker should be performed if the treatment population or digital therapeutic change substantially. And when applying these novel biomarkers, we must appreciate that unknown sources of bias may exist, so that we avoid over-reliance on such biomarkers.

**Future Work**

There are three areas of work that will extend this initial phase of digital biomarker development. As research expands into larger clinical trials, it will enable the revalidation and to some degree reconstruction of these biomarkers using larger training datasets, creating even more robust biomarkers. A larger dataset enables inclusion of other potentially predictive variables, such as demographics, sociomarkers, or omics data, and splitting of the dataset into a training and test set to minimize overfitting.

Second, it will be important to study whether the intended effects (e.g., the improvement of treatment outcomes) are actualized when these biomarkers are put into practice and to observe whether there are any unintended consequences. Alongside empirical research, software usability testing must insure that the practical application of biomarkers is interpreted by both patients and clinicians in the intended manner.

Third, similar machine learning methods can be applied to develop digital biomarkers that predict present physiological status, for instance current blood pressure or fasting blood sugar. This will require a larger training dataset with frequent (i.e., multiple times per week) measures of the ground truth (i.e., resting blood pressure, fasting blood sugar). A similar validation strategy can be used to determine the validity of the biomarkers with and without self-monitoring

of the ground truth. Our aim is to develop cuff-less blood pressure and stick-less blood glucose biomarkers that would allow for more continuous patient care at a lower burden to patients and the health system. Our hypothesis is that these biomarkers will significantly reduce clinical inertia, enhance behavioral therapy delivery, and further empower patients and providers, meaningfully increasing treatment outcomes at both the patient and population level.

## CONCLUSIONS

Machine learning can be used to transform data from a digital therapeutic into actionable digital biomarkers. In this paper, we present a successful proof of concept for a biomarker that utilizes 28 days of patient-generated data to predict a clinically meaningful response to digitally-delivered behavioral therapy. Many practical and ethical considerations arise in the development of digital biomarkers. Applying conventional clinical thinking to these novel computational processes provides the basis to identify and resolve these considerations. There is great potential to design digital biomarkers to enhance the delivery of medical care and improve treatment outcomes.

## AUTHOR CONTRIBUTONS

NLG, MAB, JC and SD prepared the data and conducted the analysis. All authors contributed to the conceptualization of the project, the interpretations of results and the writing of all sections of the paper.

## COMPETING INTERESTS

All authors have completed the ICMJE uniform disclosure form at www.icmje.org/coi_disclosure.pdf and declare: authors NG, KE, KA, DK, and MB are employees and equity owners of Better Therapeutics, LLC, authors DE, JC, and SD are independent paid scientific consultants of Better Therapeutics and JC was provided the raw data to perform all machine learning methods independently; no other relationships or activities that could appear to have influenced the submitted work.

## ETHICAL APPROVAL

Quorum Institutional Review Board, Seattle, Washington.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

**PROVENCE AND REVIEW**

Not commissioned; externally peer reviewed.

**DATA SHARING STATEMENT**

Data used for the development of biomarkers and predictive models presented here are available upon reasonable request.

**REFERENCES**

[1]     Benjamin EJ, Muntner P, Alonso A, et al. Heart disease and stroke statistics—2019 update: a report from the American Heart Association. *Circulation* 2019; 139: e56–e66. doi: 10.1161/CIR.0000000000000659

[2]     National Center for Chronic Disease Prevention and Health Promotion. *The power of prevention: chronic disease...the public health challenge of the 21st century.* CS201478, United States Department of Health and Human Services, https://www.cdc.gov/chronicdisease/pdf/2009-Power-of-Prevention.pdf (2009, accessed 27 February 2019).

[3]     Ford ES, Bergmann MM, Kröger J, et al. Healthy living is the best revenge. *Arch Intern Med* 2009; 169: 9. doi:10.1001/archinternmed.2009.237

[4]     Turner RM, Ma Q, Lorig K, et al. Evaluation of a diabetes self-management program: claims analysis on comorbid illnesses, health care utilization, and cost. *JMIR* 2018; 20: e207. doi:10.2196/jmir.9225

[5]     Kvedar JC, Fogel AL, Elenko E, et al. Digital medicine's march on chronic disease. *Nat biotechnol* 2016; 34: 239–246. doi:10.1038/nbt.3495

[6]     Charpentier G, Benhamou P-Y, Dardari D, et al. The Diabeo software enabling individualized insulin dose adjustments combined with telemedicine support improves HbA1c in poorly controlled type 1 diabetic patients. *Diabetes Care* 2011; 34: 533–539. doi:10.2337/dc10-1259

[7]     Milani RV, Lavie CJ, Bober RM, et al. Improving hypertension control and patient engagement using digital tools. *Am J Med*; 130: 14–20. doi: 10.1016/j.amjmed.2016.07.029

[8]     Quinn CC, Clough SS, Minor JM, et al. WellDoc mobile diabetes management randomized controlled trial: change in clinical and behavioral outcomes and patient and physician satisfaction. *Diabetes Technol Ther* 2008; 10: 160–168. doi: 10.1089/dia.2008.0283

[9]     Berman MA, Guthrie NL, Edwards KL, et al. Change in glycemic control with use of a digital therapeutic in adults with type 2 diabetes: Cohort Study. *JMIR Diabetes* 2018; 3: e4. doi: 10.2196/diabetes.9591

[10]    Coravos A, Khozin S, Mandl KD. Developing and adopting safe and effective digital biomarkers to improve patient outcomes. *NPJ Digit Med*; 2. Published Online First: December 2019. DOI: 10.1038/s41746-019-0090-4.

[11]    Meister S, Deiters W, Becker S. Digital health and digital biomarkers – enabling value chains on health data. *Current Directions in Biomedical Engineering*; 2. Published Online First: 1 January 2016. doi: 10.1515/cdbme-2016-0128.

[12]    Wright J, Regele O, Kourtis L, et al. Evolution of the digital biomarker ecosystem. *Digital Medicine* 2017; 3: 154. doi: 10.4103/digm.digm_35_17

[13]    Fritz BA, Chen Y, Murray-Torres TM, et al. Using machine learning techniques to develop forecasting algorithms for postoperative complications: protocol for a retrospective study. *BMJ Open* 2018; 8: e020124.

[14]    Westerman K, Reaver A, Roy C, et al. Longitudinal analysis of biomarker data from a personalized nutrition platform in healthy subjects. *Sci Rep*; 8. Published Online First: December 2018. doi: 10.1038/s41598-018-33008-7.

[15]    Minich DM, Bland JS. Personalized lifestyle medicine: relevance for nutrition and lifestyle recommendations. *The Scientific World Journal* 2013; 2013: 1–14. doi: 10.1155/2013/129841

[16]    Darcy AM, Louie AK, Roberts LW. Machine learning and the profession of medicine. *JAMA* 2016; 315: 551–552.

[17]    Sun D, Liu J, Xiao L, et al. Recent development of risk-prediction models for incident hypertension: An updated systematic review. *PLOS ONE* 2017; 12: e0187240.

[18]    Thornton RLJ, Glover CM, Cené CW, et al. Evaluating strategies for reducing health disparities by addressing the social determinants of health. *Health Aff* 2016; 35: 1416–1423. doi: 10.1377/hlthaff.2015.1357

[19]    Egger G, Dixon J. Beyond obesity and lifestyle: A review of 21st century chronic disease determinants. *BioMed Res Int* 2014; 2014: 1–12. doi: 10.1155/2014/731685

[20]    Gastil R. The determinants of human behavior. *Am Anthropol* 1961; 63: 1281–1291. http://www.jstor.org/stable/666861 (Accessed 5 Mar 2019).

[21]    Szyf M, McGowan P, Meaney MJ. The social environment and the epigenome. *Environ Mol Mutagen* 2008; 49: 46–60. doi: 10.1002/em.20357

[22]    Dagum P. Digital biomarkers of cognitive function. *NPJ Digit Med*; 1. Published Online First: December 2018. doi: 10.1038/s41746-018-0018-4.

[23]    Shin EK, Mahajan R, Akbilgic O, et al. Sociomarkers and biomarkers: predictive modeling in identifying pediatric asthma patients at risk of hospital revisits. *NPJ Digit Med*; 2018; 1: 50. doi:10.1038/s41746-018-0056-y

[24]    Billings J, Blunt I, Steventon A, et al. Development of a predictive model to identify inpatients at risk of re-admission within 30 days of discharge (PARR-30). *BMJ Open* 2012; 2: e001667. doi: 10.1136/bmjopen-2012-001667

[25]    Whelton PK, Carey RM, Aronow WS, et al. 2017 ACC/AHA/AAPA/ABC/ACPM/AGS/APhA/ASH/ASPC/NMA/PCNA Guideline for the prevention, detection, evaluation, and management of High Blood Pressure in Adults. *J Am Coll Cardiol* 2018; 71: e127–e248. doi: 10.1161/HYP.0000000000000065

[26]    Williams B, Mancia G, Spiering W, et al. 2018 ESC/ESH Guidelines for the management of arterial hypertension. *Eur Heart J* 2018; 39: 3021–3104. doi: 10.1093/eurheartj/ehy339

[27]    Quorum Review IRB: independent ethics review board. *Quorum Review IRB*, https://www.quorumreview.com/ (accessed 6 December 2017).

[28]    Kuhn M, Johnson K. *Applied Predictive Modeling*. 5th ed. Springer, 2016.

[29]  Ettehad D, Emdin CA, Kiran A, et al. Blood pressure lowering for prevention of cardiovascular disease and death: a systematic review and meta-analysis. *Lancet* 2016; 387: 957–967. doi: 10.1016/S0140-6736(15)01225-8

[30]  Thomopoulos C, Parati G, Zanchetti A. Effects of blood pressure lowering on outcome incidence in hypertension. 1. Overview, meta-analyses, and meta-regression analyses of randomized trials. *J Hypertens (Los Angel)* 2014; 32: 2285–2295. doi:10.1097/HJH.0000000000000447

[31]  Milman T, Joundi RA, Alotaibi NM, et al. Clinical inertia in the pharmacological management of hypertension. *Medicine (Baltimore)*; 97. Published Online First: 22 June 2018. doi: 10.1097/MD.0000000000011121.

[32]  Ogedegbe G. Barriers to optimal hypertension control. *J of Clin Hypertens (Greenwich)* 2008; 10: 644–646. doi:10.1111/j.1751-7176.2008.08329.x

[33]  Breiman L. Random forests. *Machine Learning* 2001; 45: 5–32. https://doi.org/10.1023/A:1010933404324

[34]  Liu M-X, Chen X, Chen G, et al. A Computational framework to infer human disease-associated long noncoding RNAs. *PLoS ONE* 2014; 9: e84408. doi: 10.1371/journal.pone.0084408

[35]  Doytchinova IA, Flower DR. VaxiJen: a server for prediction of protective antigens, tumour antigens and subunit vaccines. *BMC Bioinformatics* 2007; 7. doi: 10.1186/1471-2105-8-4

[36]  Moore RG, Brown AK, Miller MC, et al. The use of multiple novel tumor biomarkers for the detection of ovarian carcinoma in patients with a pelvic mass. *Gynecol Oncol* 2008; 108: 402–408. doi: 10.1016/j.ygyno.2007.10.017

[37]  Airola A, Pahikkala T, Waegeman W, et al. An experimental comparison of cross-validation techniques for estimating the area under the ROC curve. *Comput Stat & Data Anal* 2011; 55: 1828–1844. doi: 10.1016/j.csda.2010.11.018

[38]  Lundberg SM, Lee SI. Consistent feature attribution for tree ensembles. https://arxiv.org/abs/1706.06060 (2017, accessed 19 November 2018).

[39]  Shaban-Nejad A, Michalowski M, Buckeridge DL. Health intelligence: how artificial intelligence transforms population and personalized health. *NPJ Digit Med*; 1. Published Online First: December 2018. doi: 10.1038/s41746-018-0058-9.
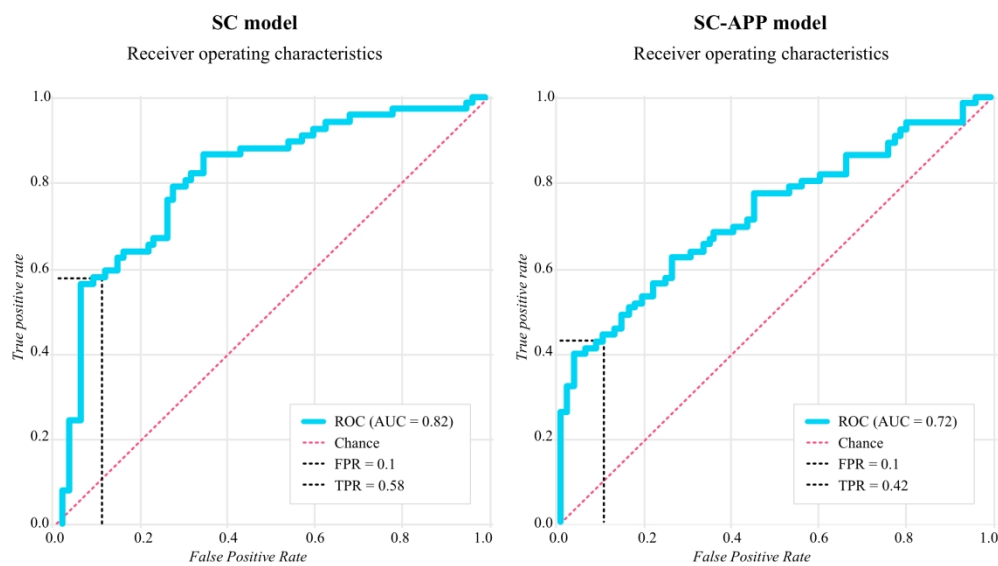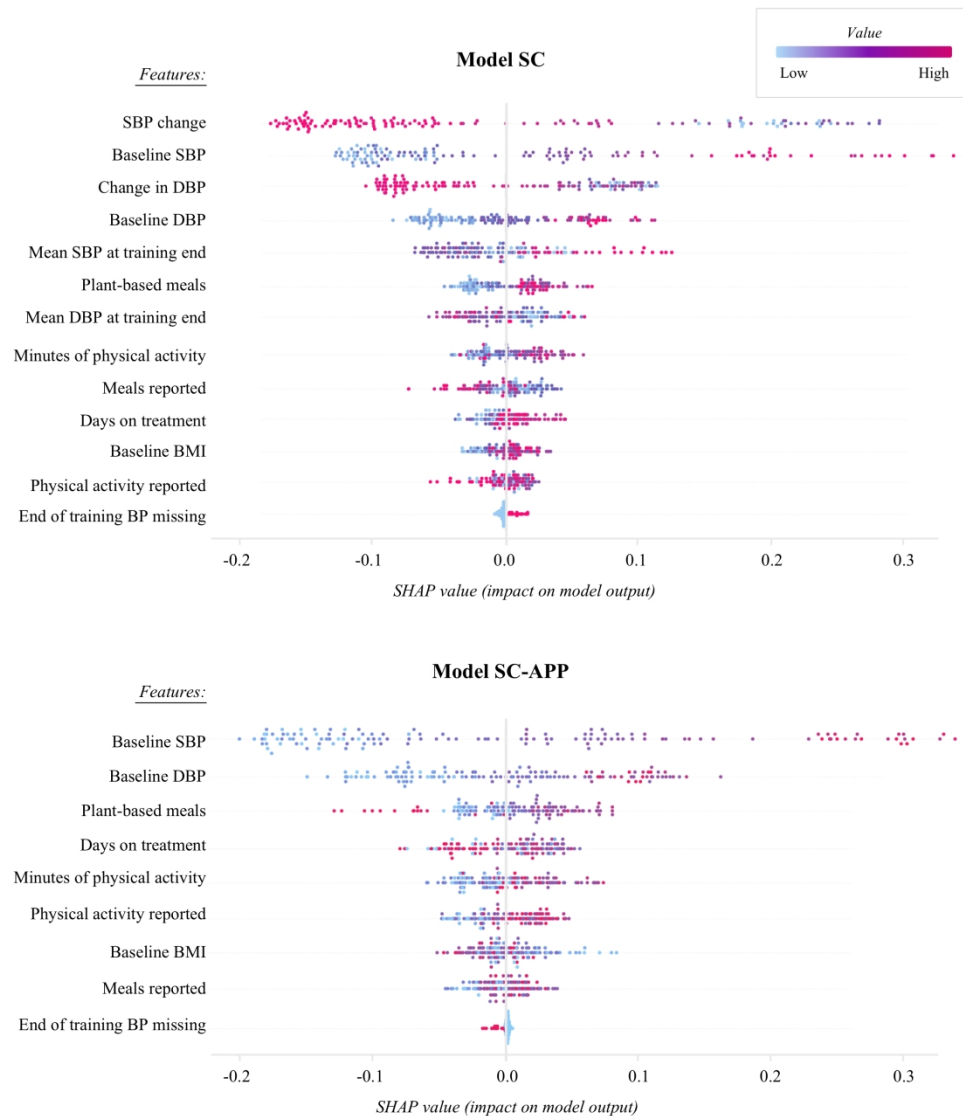
Figure 1: Receiver Operator Characteristics (ROC) curves for machine learning model predicting systolic change (SC) and a model predicting systolic change without use of ongoing blood pressure data (SC-APP).

934x608mm (72 x 72 DPI)

Figure 2: Shapley values illustrate how explanatory variables contribute to success meeting the response variable (improvement in systolic blood pressure ≥ 10 mmHg). The feature list down the y-axis is in order of contribution to the model (most to least). Each dot represents the value for one participant. SBP change and DBP change are the difference in measurements from baseline to the end of the 28-day training period.
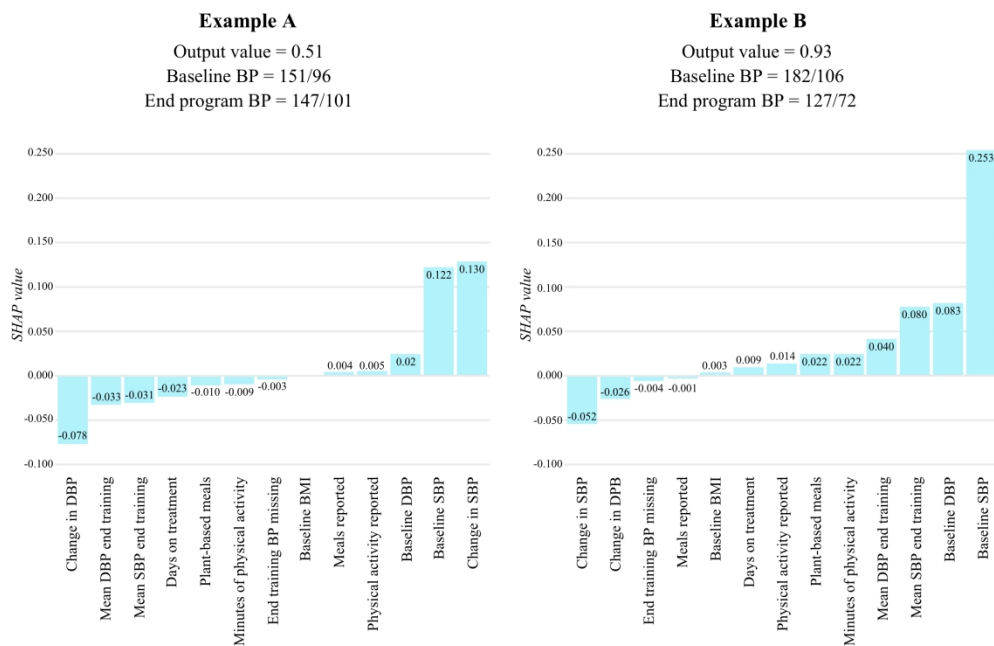
922x1058mm (72 x 72 DPI)

Figure 3: SHAP values for explanatory variables for 2 participants. The SHAP value plotted on the y-axis indicates that amount the variable positively or negatively contributes to the prediction of success (the output value). The probability threshold (output value that assigns a prediction of success) is 0.66.

934x680mm (72 x 72 DPI)

# Reporting checklist for prediction model development and validation study.

Based on the TRIPOD guidelines.

## Instructions to authors

Complete this checklist by entering the page numbers from your manuscript where readers will find each of the items listed below.

Your article may not currently address all the items on the checklist. Please modify your text to include the missing information. If you are certain that an item does not apply, please write "n/a" and provide a short explanation.

Upload your completed checklist as an extra file when you submit to a journal.

In your methods section, say that you used the TRIPOD reporting guidelines, and cite them as:

Collins GS, Reitsma JB, Altman DG, Moons KG. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): The TRIPOD statement.

|  | | Reporting Item | Page Number |
|---|---|---|---|
|  | #1 | Identify the study as developing and / or validating a multivariable prediction model, the target population, and the outcome to be predicted. | 2 |
|  | #2 | Provide a summary of objectives, study design, setting, participants, sample size, predictors, outcome, statistical analysis, results, and conclusions. | 2 |
|  | #3a | Explain the medical context (including whether diagnostic or prognostic) and rationale for developing or validating the multivariable prediction model, including references to existing models. | 3-4 |
|  | #3b | Specify the objectives, including whether the study describes the development or validation of the model or both. | 4 |
| Source of data | #4a | Describe the study design or source of data (e.g., randomized trial, cohort, or registry data), separately for the development and validation data sets, if applicable. | 4-5 |

| | | | | |
|---|---|---|---|---|
| | #4b | Specify the key study dates, including start of accrual; end of accrual; and, if applicable, end of follow-up. | | n/a |
| Participants | #5a | Specify key elements of the study setting (e.g., primary care, secondary care, general population) including number and location of centres. | | 4 |
| | #5b | Describe eligibility criteria for participants. | | 5 |
| | #5c | Give details of treatments received, if relevant | | 4 |
| Outcome | #6a | Clearly define the outcome that is predicted by the prediction model, including how and when assessed. | | 5-6 |
| | #6b | Report any actions to blind assessment of the outcome to be predicted. | | n/a |
| Predictors | #7a | Clearly define all predictors used in developing or validating the multivariable prediction model, including how and when they were measured | | 5-7 |
| | #7b | Report any actions to blind assessment of predictors for the outcome and other predictors. | | n/a |
| Sample size | #8 | Explain how the study size was arrived at. | | 5 |
| Missing data | #9 | Describe how missing data were handled (e.g., complete-case analysis, single imputation, multiple imputation) with details of any imputation method. | | n/a |
| Statistical analysis methods | #10a | If you are developing a prediction model describe how predictors were handled in the analyses. | | 5-7 |
| | #10b | If you are developing a prediction model, specify type of model, all model-building procedures (including any predictor selection), and method for internal validation. | | 6-8 |
| | #10c | If you are validating a prediction model, describe how the predictions were calculated. | | 7-8 |
| | #10d | Specify all measures used to assess model performance and, if relevant, to compare multiple models. | | 7-8 |
| | #10e | If you are validating a prediction model, describe any model updating (e.g., recalibration) arising from the validation, if done | | 7-8 |

| | | | | |
|---|---|---|---|---|
| Risk groups | #11 | Provide details on how risk groups were created, if done. | | n/a |
| Development vs. validation | #12 | For validation, identify any differences from the development data in setting, eligibility criteria, outcome, and predictors. | | 8 |
| Participants | #13a | Describe the flow of participants through the study, including the number of participants with and without the outcome and, if applicable, a summary of the follow-up time. A diagram may be helpful. | | 9 |
| | #13b | Describe the characteristics of the participants (basic demographics, clinical features, available predictors), including the number of participants with missing data for predictors and outcome. | | 9-10 |
| | #13c | For validation, show a comparison with the development data of the distribution of important variables (demographics, predictors and outcome). | | 9-10 |
| Model development | #14a | If developing a model, specify the number of participants and outcome events in each analysis. | | 9-10 |
| | #14b | If developing a model, report the unadjusted association, if calculated between each candidate predictor and outcome. | | 10-11 |
| Model specification | #15a | If developing a model, present the full prediction model to allow predictions for individuals (i.e., all regression coefficients, and model intercept or baseline survival at a given time point). | | 10-11 |
| | #15b | If developing a prediction model, explain how to the use it. | | 10-11 |
| Model performance | #16 | Report performance measures (with CIs) for the prediction model. | | 10-11 |
| Model-updating | #17 | If validating a model, report the results from any model updating, if done (i.e., model specification, model performance). | | n/a |
| Limitations | #18 | Discuss any limitations of the study (such as nonrepresentative sample, few events per predictor, missing data). | | 12-13 |
| Interpretation | #19a | For validation, discuss the results with reference to performance in the development data, and any other validation data | | 11-12 |
| | #19b | Give an overall interpretation of the results, considering objectives, limitations, results from similar studies, and other relevant evidence. | | 14 |

| | | | |
|---|---|---|---|
| Implications | #20 | Discuss the potential clinical use of the model and implications for future research | 13-14 |
| Supplementary information | #21 | Provide information about the availability of supplementary resources, such as study protocol, Web calculator, and data sets. | 14-15 |
| Funding | #22 | Give the source of funding and the role of the funders for the present study. | 14 |

The TRIPOD checklist is distributed under the terms of the Creative Commons Attribution License CC-BY. This checklist was completed on 27. March 2019 using https://www.goodreports.org/, a tool made by the EQUATOR Network in collaboration with Penelope.ai

# Emergence of digital biomarkers to predict and modify treatment efficacy: a machine learning study

SCHOLARONE™
Manuscripts

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

# Emergence of digital biomarkers to predict and modify treatment efficacy: a machine learning study

Nicole L Guthrie, MS
Better Therapeutics LLC
San Francisco, California, United States

Jason Carpenter, MS
Manifold, Inc.
Oakland, CA 94607

Katherine L Edwards, FNP-C
Better Therapeutics LLC
San Francisco, California

Kevin J Appelbaum, BSE
Better Therapeutics LLC
San Francisco, California

Sourav Dey, PhD
Manifold, Inc.
Oakland, CA

David M. Eisenberg, MD
Department of Nutrition, Harvard T.H. Chan School of Public Health
Boston, Massachusetts, United States

David L. Katz, MD, MPH
Yale University Prevention Research Center, Griffin Hospital
Derby, CT, United States

Mark A. Berman, MD [corresponding author]
Better Therapeutics LLC
445 Bush Street, Suite 300
San Francisco, California 94108
617-877-0327, mark@bettertherapeutics.io

Word count: 5089

**ABSTRACT**

**Objectives:** Development of digital biomarkers to predict treatment response to a digital behavioral intervention.

**Design:** Machine learning using random forest classifiers on data generated through the use of a digital therapeutic which delivers behavioral therapy to treat cardiometabolic disease. Data from 13 explanatory variables (biometric and engagement in nature) generated in the first 28 days of a 12-week intervention were used to train models. Two levels of response to treatment were predicted: 1) systolic change $\geq 10$ mmHg (SC model), and 2) shift down to a blood pressure category of elevated or better (ER model). Models were validated using leave-one-out cross-validation and evaluated using area under the curve receiver operating characteristics (AUROC) and specificity-sensitivity. Ability to predict treatment response with a subset of 9 variables, including app use and baseline blood pressure was also tested (models SC-APP and ER-APP).

**Setting:** Data generated through ad libitum use of a digital therapeutic in the United States.

**Participants:** De-identified data from 135 adults with a starting blood pressure $\geq 130 / 80$, who tracked blood pressure for at least 7 weeks using the digital therapeutic.

**Results:** The SC model had an AUROC of .82 and a sensitivity of 58% at a specificity of 90%. The ER model had an AUROC of .69 and a sensitivity of 32% at a specificity at 91%. Dropping explanatory variables related to blood pressure resulted in an AUROC of .72 with a sensitivity of 42% at a specificity of 90% for the SC-APP model and an AUROC of .53 for the ER-APP model.

**Conclusions:** Machine learning was used to transform data from a digital therapeutic into digital biomarkers that predicted treatment response in individual participants. Digital biomarkers have potential to improve treatment outcomes in a digital behavioral intervention.

**ARTICLE SUMMARY**

**Strengths and Limitations**

- Proof of concept biomarkers demonstrated good power to predict treatment outcomes despite the small size of the training dataset.
- Use of additional explanatory variables to develop the biomarkers may enhance the accuracy of predictions.
- Generalizability of the biomarkers is unknown and may be limited by the demographics of the training dataset.

**INTRODUCTION**

Modifiable behaviors are responsible for 70% or more of all cardiometabolic diseases.[1–3] Health systems are ill equipped to address the current epidemic of cardiometabolic diseases and, in particular, lack widely available behavioral therapies to address the common root causes of these conditions. Digital therapeutics, software designed to encourage changes in behaviors which treat disease, offer a means to deliver behavioral therapy to large populations and

preliminary studies demonstrate their potential as a cost-effective treatment for cardiometabolic disease.[4–7]

Compared to pharmacotherapy, digital therapeutics offer potential advantages such as ease of access, ease of use, fewer side-effects and cost-effectiveness.[4, 8–10] Digital therapeutics also generate readily accessible patient data without requiring an office or lab visit. The data generated are voluminous and vary in both type and quality. These can include remotely sensed measures of physiology (e.g., blood pressure, blood glucose, heart rate variability), behavioral data (e.g., about eating, moving, thinking), medication adherence, as well as engagement parameters of unknown significance (e.g., app use, geographic location, circadian patterns of use).

The best use of these data remains an open question. Feeding data directly into EMRs is of limited utility to providers or patients. Whereas, transforming the data into markers of disease status, termed digital biomarkers[11] could provide clinically actionable insights with or without conventional biometric data.[12–14] Digital biomarkers afford a pragmatic approach to remotely monitor patients and intervene on a continuous rather than episodic basis. Greatly expanding opportunities to intervene means that patients have greater access to personalized care, which could improve treatment outcomes.[15, 16]

Machine learning, a type of artificial intelligence (AI) used to make predictions with large and complex datasets, offers a novel approach for creating digital biomarkers. The exponential growth of smartphone use in the United States and advancing interoperability standards allow for digital biomarkers to be compiled across diverse populations and data sources. As a result, the opportunity to advance digital biomarker methodologies has never been greater.[17, 18]

Machine learning is particularly valuable when there is ambiguity about what variables, or to what extent a set of variables predict an outcome of interest. Such ambiguity is inherent to behavioral interventions, like those used in the treatment of cardiometabolic disease. Human behavior results from a complex interaction between the anthropogenic features of our living environment, genetic and epigenetic determinants of behavior, neurobiology, social influences and to some degree chance events.[19–22] While clinical experience and the scientific literature can identify many variables that influence behavior, the interplay of these variables in a given individual and their environment is difficult for a clinician or patient to discern. This ambiguity limits enthusiasm for behavioral interventions because it makes it difficult for clinicians and patients to rely on behavioral therapy to achieve a predictable level of therapeutic response.

Digital biomarkers can reduce ambiguity by predicting current and forecasting future disease status during the course of treatment.[23–25] Digital biomarkers that serve as markers of current disease status allow for tailoring or adjusting treatment between clinic visits (e.g., when a patient

is not doing as well as expected). Similarly, markers of future disease status enable preemptive action, such as adding or subtracting additional treatments, or taking preventive steps to avoid complications of the disease.[14, 25]

In this paper, we present our ongoing work to develop digital biomarkers and aim to illustrate their utility in digitally-delivered behavioral interventions to both the patient and prescribing clinician. We describe the analytic process to show that the development of digital biomarkers requires a hypothesis-driven approach, particularly when datasets are small. Finally, using actual examples of digital biomarkers intended to predict blood pressure status, we discuss the practical and ethical considerations involved in both developing and applying digital biomarkers using machine learning.

## METHODS
### Digital Therapeutic
The digital biomarkers discussed in this paper were generated using data from a digital therapeutic created by Better Therapeutics LLC (San Francisco, CA), a developer of prescription digital therapeutics for the treatment of cardiometabolic diseases.

The digital therapeutic integrates a mobile medical application ("app") that delivers behavioral therapy with the support of a remote multidisciplinary care team. The app delivers a personalized behavior change intervention, including tools for goal setting, skill building, self-monitoring, biometric tracking and behavioral feedback designed to provide cognitive training and support the participant's daily efforts to improve overall cardiometabolic disease status. The app facilitates the adoption of evidenced-based behavioral strategies, such as planning and self-monitoring, to increase physical activity and consumption of vegetables, whole grains, fruits, nuts, seeds, beans and other legumes.[26, 27]

### Participants
Participants were over 18 years of age who self-identified with a diagnosis of hypertension, type 2 diabetes and/or hyperlipidemia. Those reporting to have hypertension at enrollment were offered a free Omron 7 Series Upper Arm Blood Pressure Monitor (Omron Healthcare Inc, Kyoto, Japan) to use during the intervention and to keep at its conclusion. This 12-week intervention was available to all participants at no cost.

This analysis of existing data from participants using the digital therapeutic was approved and overseen by Quorum Review Institutional Review Board[28] and a waiver of informed consent was granted for this retrospective analysis.

### Constraining the Training Dataset

The clinical intent of this biomarker exploration was to generate digital biomarkers that could improve treatment outcomes in future participants with hypertension for whom blood pressure was not optimally controlled despite the use of pharmacological treatment. We first identified prior participants that met the following criteria: baseline blood pressure was 1) at or above the cutoff for stage I hypertension (systolic $\geq$ 130 or diastolic $\geq$ 80)[26] and, 2) recorded no more than 2 weeks prior to or 2 weeks after the start of the intervention. The baseline value was calculated by taking an average of all values reported in a 6-day interval defined as starting with the date of the first blood pressure value reported and all values reported in the following 5 days.

The development of a predictive model using a training dataset requires all participants to have known outcomes. This is referred to as the "ground truth" in machine learning. In this case, the ground truth was change in blood pressure from baseline. The follow-up blood pressure values were calculated by taking an average over a 6-day interval ending with the last blood pressure reported and all values in the previous 5 days. Of the participants identified, we only included those with a follow-up blood pressure value in weeks 7 to 14 of the intervention.

The training set also needs to be a valid representation of future experience to make sure the data sources present are sufficiently representative of those that will be collected in future participants. This is an important consideration because of the evolving nature of digital therapeutics. As the therapeutic is modified over time, the data types collected in earlier versions of the app could be different from those collected in later versions. Therefore, for the purpose of developing biomarkers that predict blood pressure status, we further constrained the training dataset to include participants who used versions of the app which enabled automatic blood pressure tracking.

**Choosing Response Variables and Training Window**
The intent of machine learning is to train an algorithm that can predict a specific outcome, termed the "response variable." The response variable must be both clinically relevant and sufficiently, but not universally, prevalent in the population of interest. For example, if the outcome is "any degree of improvement", but "any degree of improvement" occurs in > 95% of participants in the training dataset, then a predictive model may appear to work well, but could actually be invalid.[29] Furthermore, "any improvement" is arguably less clinically meaningful than "a specific degree of improvement".

To address these concerns, we chose response variables that were well distributed in the training dataset, as detailed in the results section, and were also clinically meaningful in the management of hypertension. Specifically, we chose to predict: 1) a systolic blood pressure improvement of 10 mmHg or higher near the end of a 12-week intervention period (defined as 7 to 14 weeks after start), because 10 mmHg has been well demonstrated to be clinically meaningful[30, 31] and that degree of change is near the mean for participants in the training dataset (66/135 participants);

and 2) a reduction in blood pressure to the elevated range or below (systolic ≤ 130), because this level of blood pressure control would signal a clinician to stop adding additional pharmaceuticals or consider reducing or deprescribing pharmaceuticals (48/135 participants).[26]

For a digital biomarker that predicts blood pressure status to be clinically useful, it needs to compute with data collected in a short period of time, and in less time than typically occurs between clinic visits. This means the biomarker could be used to intervene between office visits and could play a role in addressing clinical inertia that limits primary care providers ability to optimize blood pressure control in their patients.[32, 33] To demonstrate proof of concept, we chose a 28-day training interval, meaning that we trained machine learning models on the first 28 days of patient data to evaluate whether it could predict blood pressure change in weeks 7 to 14. We hypothesized that data collected within this short training window could sufficiently represent changing behavioral patterns and treatment response, so as to predict future blood pressure status.

**Choosing Modeling Techniques and Explanatory Variables**
There are many valid methodologies for machine learning. These methods can be categorized by the type of learning involved - supervised or unsupervised. While a discussion on the merits and nature of each type is beyond the scope of this paper, it is important to select modeling techniques appropriate for the size of the dataset and nature of response variables chosen. For the biomarkers presented herein, we utilized random forest models, which is a form of supervised classification learning known to reduce overfitting in small datasets.[34] The models are supervised because the ground truth (i.e., actual blood pressure change) for each participant in the training dataset is labeled. The models are attempting to learn whether the classification was or was not achieved; for example, will a participant achieve a > 10 mmHg blood pressure reduction, or not?

In addition to classification labels, random forest models must be trained on a set of explanatory variables. For a small training dataset, each model must have a limited number of variables so as to avoid excess noise and overfitting that can lead to reduced generalizability. These explanatory variables must be selected by hypothesis. For example, we hypothesized that baseline blood pressure and achievement of behavioral goals would influence the degree of blood pressure change observed and used these as explanatory variables. In a large dataset, feature engineering can be used to identify the most predictive explanatory variables.

For each biomarker model, denoted SC (systolic change of 10 mmHg or more) or ER (elevated range achieved), we used 13 explanatory variables, which can be categorized as engagement or biometric variables. Engagement variables are counts of actions related to the use of the digital therapeutic, including count of all meals reported, plant-based meals reported, physical activity reported and length of exposure to the intervention. Biometric variables included baseline

systolic, baseline diastolic, mean systolic and diastolic at training window end, initial systolic and diastolic change (end training mean - baseline), minutes of physical activity, and baseline Body Mass Index (BMI).

We trained each biomarker model on data generated during a 28-day training window, starting with participant day 0 (sign up) up to day 28 (a third of the way through the studied 12-week intervention period). We utilized hyperparameter optimization [35] to minimize overfitting and to achieve the maximal leave-one-out cross-validated area under the receiver operating characteristic curve for both models.

**Model Validation**
Performance of each biomarker model was assessed using leave-one-out cross-validation, which is a common and valid technique for use in samples of this size.[36–38] This was done by training each model on N-1 samples of the data and then making a prediction on that one sample that was left out, producing an "out of sample' prediction for all N samples. The N predictions were pooled to generate the classification variables of the receiver operator characteristic curve (ROC), the area under the curve of the ROC (AUROC) and a confusion matrix of true versus predicted values.[39]

For each biomarker model, the ROC curve illustrates predictive ability of the response variable (in this case systolic change of 10 mmHg or move to a range of elevated BP) at different thresholds of discrimination. At each prediction discrimination threshold, the ROC displays the false positive rate (FPR) against the true positive rate (TPR). The FPR is the ratio of truly negative events categorized as positive (FP) to the total number of actual negative events (N). Specificity or true negative rate of a model is calculated as 1 - FPR and is an indication of how well a model does in correctly identifying those who do not achieve a successful outcome, as defined by the response variable.

Since the intended application of these biomarkers is analogous to a diagnostic test, which are traditionally evaluated based on their specificity, we evaluated model performance at a specificity of 90% (FPR = .10). A low FPR minimizes the number of participants who the model would predict to achieve a healthier state who actually will not. In turn, this minimizes the number of participants who might be erroneously taken off blood pressure medicine as a result of an erroneous prediction. It is less critical to avoid labeling participants who had achieved a healthier state as though they had not. This is why we choose a discrimination threshold with a low FPR, and then evaluate the TPR at that point.

As a further validation step, we examined the performance of each biomarker by excluding the four explanatory variables that capture blood pressure change in the training window. While we hypothesized that these models would perform less well, they serve to test the concept that a

digital biomarker that predicts blood pressure status can be generated without using ongoing blood pressure data. These validation models using only app engagement and other biometric variables, are denoted SC-APP and ER-APP.

## Making Use of Explainable Artificial Intelligence (AI)

Digital biomarkers that are generated using machine learning do not need to be viewed as a black box. Instead, explainable AI techniques are available that can provide more granular details about the explanatory variables that influenced the prediction made. Explainable AI can afford both individual participant and population level insights.

We used the Tree Shapley Additive Explanation (SHAP) algorithm on the random forest models[40] to generate more interpretable predictions at the participant level. The SHAP algorithm assigns each explanatory variable an importance value for each prediction. Using SHAP on a machine learning model is analogous to coefficient analysis in classical regression. Similar to coefficient analysis, it can be used to determine the relative importance of explanatory variables in addition to determining which explanatory variables drove a particular prediction. Predictions start at a base value that is the expectation of the response variable. For binary classification models, this is defined by the proportion of outcomes by class (e.g., the proportion of participants who successfully reduced their blood pressure). Then SHAP values attribute to each explanatory variable the change in expected model prediction given the addition of that explanatory variable. This provides insight into how much each explanatory variable positively or negatively impacts the prediction made for each participant. The final prediction probability of whether the participant will achieve the response variable is the sum of the base value and all of the explanatory variable attributions.

Individual SHAP values can be used to provide specific behavioral feedback to participants, with the intent of motivating a change in behavioral pattern that may improve treatment outcomes. In particular, explanatory variables that are theoretically modifiable (such as minutes of exercise, or number of plant-based meals consumed) can be displayed to motivate changes, whereas fixed explanatory variables (such as baseline values) can be displayed to provide context. SHAP values for all participants can also be plotted to reveal the overall ranking of variables in the population studied. These variable rank lists can then inform hypotheses about how to further improve the design of the digital therapeutic to optimize clinical outcomes.

All machine learning model development was done using open-source packages in Python. The packages include but are not limited to Scikit-Learn, SHAP, Pandas, and Numpy.

## Patient and Public Involvement

We did not directly involve patients or public (PPI) in the current study. However, the digital therapeutic that is the source of the database used in this study was developed with PPI input over the product development lifecycle.

## RESULTS

### Dataset

The training dataset contained 135 participants who met the inclusion criteria. The mean age was 54.9 years (95% CI 53.5, 56.3), mean baseline BMI was 34.5 (95% CI 33.1, 35.8) and 83% (112/135) were female. Based on the 2017 ACC/AHA definition,[26] half of the participants (68/135) had stage 1 hypertension at baseline, with the other half (67/135) having stage II hypertension at baseline. Of those with stage 1 hypertension at baseline, 51.5% (35/68) had isolated diastolic hypertension (i.e., diastolic BP 80-90 mmHg). Of those with stage II hypertension at baseline, 14.9% (10/67) had isolated diastolic hypertension (i.e., diastolic BP ≥ 89 mmHg). On average, participants contributed 3 blood pressure readings to the baseline value (95% CI 2.5, 3.4) and 2.5 readings to the end-intervention value (95% CI 2.2, 2.9). Baseline characteristics are listed in Table 1.

Over the intervention period examined, systolic blood pressure changed by -12.7 mmHg (95% CI -14.8, -9.6), diastolic blood pressure changed by -7.1 mmHg (95% CI -9.0, -5.2) and the mean duration of days between baseline and most recent value for blood pressure was 79.3 days (95% CI 76.8, 81.9). Of all participants, 35.6% (48/135) shifted to a blood pressure range below stage I (<130/80). A shift to a normal range was seen in 16.4% (11/67) of those starting with stage II hypertension and 29.4% (20/68) of those starting with stage I hypertension.

**Table 1**. Participant characteristics at baseline

| Participant Characteristics | Total (n = 135) | Stage I BP (n = 68) | Stage II BP (n = 67) |
|---|---|---|---|
| | | | |
| Age (years), mean (95% CI) | 54.9 (53.5 to 56.3) | 55.7 (53.7 to 57.7) | 54.2 (52.1 to 56.2) |
| Body mass index (kg/m$^2$), mean (95% CI) | 34.5 (33.1 to 35.8) | 33.8 (31.9 to 35.6) | 35.2 (33.2 to 37.2) |
| Female, n (%) | 112 (83) | 54 (79.4) | 58 (86.6) |
| Systolic BP (mmHg), mean (95% CI) | 138.9 (136.2 to 141.6) | 127.9 (126.1 to 129.7) | 150.0 (146.6 to 153.5) |

| | | | |
|---|---|---|---|
| Diastolic BP (mmHg), mean (95% CI) | 87.8 (86.1 to 89.4) | 82.3 (81.1 to 83.6) | 93.3 (90.7 to 95.8) |
| Isolated diastolic hypertension, n (%) | 45 (33.3) | 35 (51.5) | 10 (14.9) |
| BP medications (count), mean (95% CI) | 1.3 (1.1 to 1.5) | 1.2 (1.0 to 1.4) | 1.4 (1.1 to 1.7) |

**Predictive Models**

The random forest classifier achieved optimal performance with 100 trees and a minimum of 3 samples per leaf node for the SC model. For the ER model optimal performance was achieved with 400 trees and a minimum of 5 samples per leaf nodes.

Biomarker models were assessed at the operating point on each ROC that was as close as possible to a FPR of 10%. The SC model (predicting a systolic change of 10 points) was assessed at a FPR of 10%, which means that 10% of participants who didn't achieve a reduction in systolic blood pressure of 10 mmHg were labeled as though they had. Evaluating the model at 10% FPR, we were able to achieve a TPR of 58%. This means that 58% of participants who achieved a reduction in systolic blood pressure of 10 mmHg were labeled correctly. The AUROC was .82, model specificity (1 - FPR) was 90%, sensitivity (TPR) was 58% and accuracy ((TP+TN)/n) was 74%. In the SC-APP model, where variables related to changes in blood pressure were removed, the AUROC was .72 and at a FPR of 10% (specificity of 90%), the TPR was 42%. The resultant receiver operator curves for these 2 models can be seen in Figure 1.

The biomarker models exploring the ability to predict a shift down to a blood pressure range of elevated or better (ER and ER-APP) also demonstrated predictive capacity, but less so than the SC models. For the ER model, the AUROC was .69 and at a FPR of 9% the TPR was 32%. When the BP change variables were removed, in the ER-APP model, the prediction ability was only slightly above chance (AUC = .53, TPR 26% at FPR of 12%).

Plots of the Tree SHAP algorithm results for the SC and SC-APP models are shown in Figure 2. Explanatory variables on the y-axis are ordered from most to least predictive based on their average absolute contribution to the response variable. Each dot represents the SHAP value of that variable for one participant and the placement of the dots on the x-axis indicate if the contribution was subtractive or additive for a specific participant. The color of the dot is indicative of the value for that variable, with highly positive values displayed as red and low or negative values showing up as light blue. The plot for the SC model reveals that explanatory

variables related to blood pressure were top contributors to the prediction. For example, the distribution of dots across the x-axis for the first variable listed shows that improvements in systolic BP early in the intervention, as seen by the blue and purple dots to the right of 0 on the x-axis, contributed positively to the prediction that a participant would succeed. Behavioral variables also had predictive power. For example, a high count of physical activity minutes and plant-based meals reported positively contributed to a prediction of success for most participants.

Shapley values can be aggregated and illustrated for every participant. A plot of the SHAP values helps to visualize which variables contributed most to a low or high prediction of success for an individual participant. In figure 3, we display the SHAP values for two participants, one with a lower than expected probability of success (example A), and one with a higher than expected probability of success (example B). In example A, the participant experienced a large improvement in their systolic blood pressure in weeks 3 and 4 (-14 mmHg), yet is given a low probability of sustaining this improvement at the end of the intervention period. This surprisingly low probability is explained by the SHAP values, which reveal low counts for several behavioral explanatory variables, such as the number of plant-based meals and minutes of physical activity reported. This data can be automatically translated into a simple explanation to the participant, that their probability for sustaining meaningful change could be higher if they made incremental improvements in their meal and activity pattern. Furthermore, the exact number of additional meals and activity minutes to accrue per week to sufficiently increase the probability of success could be calculated, to motivate the participant and to give meaning to the additional efforts prescribed.

In example B, the participant has evidenced no improvement in systolic blood pressure at the end of week 4, yet they are predicted to meaningfully improve blood pressure by the of the treatment period. This unexpected prediction is explained by the SHAP values, which show that the combined impact of their baseline blood pressure and behavioral explanatory variables suggests a high likelihood of success. This data can be automatically translated to provide timely encouragement to the participant to maintain or advance their behavioral changes even though their blood pressure has not yet responded.

**DISCUSSION**

In this paper, we present a proof of concept that digital biomarkers can be developed using even a small training dataset from a digital therapeutic. We demonstrated that 28 days of data can be transformed, using machine learning, into a digital biomarker that predicts the degree of treatment response, in this case whether a meaningful drop in blood pressure will occur at the end of the treatment period.

There are many ways to use such a biomarker in practice to tailor behavioral treatment and improve outcomes. For a patient, these biomarkers can be used as a continuous form of treatment

feedback and behavioral reinforcement. The probability of a significant treatment response can be translated into a treatment score, much like a credit score. Since this score could be recalculated with every new engagement recorded in the digital therapeutic, it would serve to motivate app use and reinforce healing behaviors. In addition, the biomarker output can be made more meaningful using explainable AI techniques. For example, SHAP values can be translated into a prioritized list of behavioral actions to help a patient focus their attention on efforts that are most predictive of success.

For clinicians and health systems, digital biomarkers can function as a form of automated patient monitoring. The probability of a positive treatment response can be translated into a clinical alert by setting an acceptable specificity-sensitivity threshold for each biomarker paired with a duration of time above this threshold. Like any diagnostic test, the performance characteristics of the alert should be made known to those acting upon it. Since such an alert would be intended to influence treatment decisions, for example via a clinical decision support tool, specificity-sensitivity pairs need to be evaluated from a risk-benefit perspective. For instance, how do the risks associated with false positives and false negatives compare to the benefits of identifying true positives and true negatives? To accurately weigh these risks and benefits requires us to understand the context that the biomarker and therapeutics are used. Current clinical practice, for example, is plagued by high rates of clinical inertia (i.e., a lack of timely and appropriate treatment decisions). Therefore, a higher false positive rate may be tolerated as a trade-off for an easier-to-access biomarker.

For a developer of digital therapeutics, digital biomarkers provide not just a way to personalize treatment and communicate clinical status to providers, but also a way to better understand what variables within the therapeutic are most predictive of clinical outcomes. These data can be used to guide the ongoing refinement of a digital therapeutic. When datasets are of sufficient size, the machine learning techniques used to generate digital biomarkers can also be applied to identify distinct digital phenotypes, that is, unique patterns of engagement with a behavioral intervention that represent meaningful subpopulations who share the same diagnosis. Identifying and targeting treatment to previously unknown subpopulations is thought of as meaningful step towards more personalized medicine.[16, 41]

**Limitations, Practical and Ethical Considerations**
The main limitation of the work presented here is the size of the training dataset used. It is likely that a larger dataset would improve the performance characteristics of biomarkers tested.[15] It is noteworthy, given this limitation, that one of the biomarkers (SC) had an AUC greater than 0.8. This suggests the utility of beginning digital biomarker development in early in the implementation of a digital therapeutic.

To lower the risk of prematurely taking patients off of medications, the digital biomarkers presented decreased the false positive rate (i.e., a higher specificity), which resulted in a lower true positive rate (i.e., a lower sensitivity). In this context, a lower sensitivity means that the digital biomarker will fail to identify some successful individuals and as a result they may not have their medications reduced as promptly as possible. This means that the current performance of the digital biomarker does not fully obviate the need for traditional biomarkers or in-office visits.

Other limitations include the omission of known predictors of treatment response (e.g., time since diagnosis, medication adherence or change), the reliance on a small set of explanatory variables and the inclusion of self-reported variables that may be subject to human error. Addressing these limitations may enable more accurate biomarkers.

The ethical and practical implications of applying complex, ever-changing, predictive models generated by machine learning are only beginning to be appreciated. To preempt potential misuse of digital biomarkers, we must understand the true meaning of the data these biomarkers present. For example, a predictive model can identify variables that are predictive of a given outcome. This does not mean highly predictive variables caused the outcome, nor does it mean that poorly predictive variables are not causative. Instead, the predictive strength of each variable should be treated literally as "markers." This does not preclude the automated use of explanatory variables to guide the personalization of behavioral therapy. However, since the individual level of each variable could be influenced by unknown confounding factors, and since the degree of modifiability of each variable is not known, the impact of this form of behavioral therapy must be studied.

Finally, like any biomarker, a digital biomarker is only generalizable if the training dataset is truly representative of future patients. If the training dataset is biased or overly skewed, it may produce a biomarker that underperforms at best, and is harmful at worst. To guard against this bias, re-validation of a digital biomarker should be performed if the treatment population or digital therapeutic change substantially. And when applying these novel biomarkers, we must appreciate that unknown sources of bias may exist, so that we avoid over-reliance on such biomarkers.

**Future Work**
There are three areas of work that will extend this initial phase of digital biomarker development. As research expands into larger clinical trials, it will enable the revalidation (often called external validation) and to some degree reconstruction of these biomarkers using larger training datasets, creating even more robust biomarkers. A larger dataset enables inclusion of other potentially predictive variables, such as demographics, sociomarkers, or omics data, and splitting of the dataset into a training and test set to minimize overfitting. External validation also gives

the end-users of the biomarker greater confidence that the biomarker preforms well with varied individuals, settings or time or year.[42–44]

Second, it will be important to conduct an impact analysis to study whether the intended effects (e.g., the improvement of treatment outcomes) are actualized when these biomarkers are put into practice and to observe whether there are any unintended consequences. Alongside empirical research, software usability testing must insure that the practical application of biomarkers is interpreted by both patients and clinicians in the intended manner.

Third, similar machine learning methods can be applied to develop digital biomarkers that predict present physiological status, for instance current blood pressure or fasting blood sugar. This will require a larger training dataset with frequent (i.e., multiple times per week) measures of the ground truth (i.e., resting blood pressure, fasting blood sugar). A similar validation strategy can be used to determine the validity of the biomarkers with and without self-monitoring of the ground truth. Our aim is to develop cuff-less blood pressure and stick-less blood glucose biomarkers that would allow for more continuous patient care at a lower burden to patients and the health system. Our hypothesis is that these biomarkers will significantly reduce clinical inertia, enhance behavioral therapy delivery, and further empower patients and providers, meaningfully increasing treatment outcomes at both the patient and population level.

## CONCLUSIONS

Machine learning can be used to transform data from a digital therapeutic into actionable digital biomarkers. In this paper, we present a successful proof of concept for a biomarker that utilizes 28 days of patient-generated data to predict a clinically meaningful response to digitally-delivered behavioral therapy. Many practical and ethical considerations arise in the development of digital biomarkers. Applying conventional clinical thinking to these novel computational processes provides the basis to identify and resolve these considerations. There is great potential to design digital biomarkers to enhance the delivery of medical care and improve treatment outcomes.

## ACKNOWLEDGMENTS

## CONTRIBUTORS

NLG, MAB, JC and SD prepared the data and conducted the analysis. NLG, MAB, KE, KA, JC, SD, DLK, DME contributed to the conceptualization of the project and the interpretations of results. NLG, MAB, KE, KA, JC and SD contributed to the writing of the paper.

**ETHICAL APPROVAL**

Quorum Institutional Review Board, Seattle, Washington.

**PROVENCE AND REVIEW**

Not commissioned; externally peer reviewed.

**DATA SHARING STATEMENT**

Data used for the development of biomarkers and predictive models presented here are available
upon reasonable request.

**REFERENCES**

[1]     Benjamin EJ, Muntner P, Alonso A, et al. Heart disease and stroke statistics—2019
        update: a report from the American Heart Association. *Circulation* 2019; 139: e56–e66.
        doi: 10.1161/CIR.0000000000000659

[2]     National Center for Chronic Disease Prevention and Health Promotion. *The power of
        prevention: chronic disease...the public health challenge of the 21st century*. CS201478,
        United States Department of Health and Human Services,
        https://www.cdc.gov/chronicdisease/pdf/2009-Power-of-Prevention.pdf (2009, accessed
        27 February 2019).

[3]     Ford ES, Bergmann MM, Kröger J, et al. Healthy living is the best revenge. *Arch Intern
        Med* 2009; 169: 9. doi:10.1001/archinternmed.2009.237

[4]     Turner RM, Ma Q, Lorig K, et al. Evaluation of a diabetes self-management program:
        claims analysis on comorbid illnesses, health care utilization, and cost. *JMIR* 2018; 20:
        e207. doi:10.2196/jmir.9225

[5]     Kvedar JC, Fogel AL, Elenko E, et al. Digital medicine's march on chronic disease. *Nat
        biotechnol* 2016; 34: 239–246. doi:10.1038/nbt.3495

[6]     Charpentier G, Benhamou P-Y, Dardari D, et al. The Diabeo software enabling
        individualized insulin dose adjustments combined with telemedicine support improves
        HbA1c in poorly controlled type 1 diabetic patients. *Diabetes Care* 2011; 34: 533–539.
        doi:10.2337/dc10-1259

[7]    Wang J, Cai C, Padhye N, et al. A Behavioral Lifestyle Intervention Enhanced With Multiple-Behavior Self-Monitoring Using Mobile and Connected Tools for Underserved Individuals With Type 2 Diabetes and Comorbid Overweight or Obesity: Pilot Comparative Effectiveness Trial. *JMIR mHealth and uHealth* 2018; 6: e92. doi: 10.2196/mhealth.4478

[8]    Milani RV, Lavie CJ, Bober RM, et al. Improving hypertension control and patient engagement using digital tools. *Am J Med*; 130: 14–20. doi: 10.1016/j.amjmed.2016.07.029

[9]    Quinn CC, Clough SS, Minor JM, et al. WellDoc mobile diabetes management randomized controlled trial: change in clinical and behavioral outcomes and patient and physician satisfaction. *Diabetes Technol Ther* 2008; 10: 160–168. doi: 10.1089/dia.2008.0283

[10]   Berman MA, Guthrie NL, Edwards KL, et al. Change in glycemic control with use of a digital therapeutic in adults with type 2 diabetes: Cohort Study. *JMIR Diabetes* 2018; 3: e4. doi: 10.2196/diabetes.9591

[11]   Coravos A, Khozin S, Mandl KD. Developing and adopting safe and effective digital biomarkers to improve patient outcomes. *NPJ Digit Med*; 2. Published Online First: December 2019. DOI: 10.1038/s41746-019-0090-4.

[12]   Meister S, Deiters W, Becker S. Digital health and digital biomarkers – enabling value chains on health data. *Current Directions in Biomedical Engineering*; 2. Published Online First: 1 January 2016. doi: 10.1515/cdbme-2016-0128.

[13]   Wright J, Regele O, Kourtis L, et al. Evolution of the digital biomarker ecosystem. *Digital Medicine* 2017; 3: 154. doi: 10.4103/digm.digm_35_17

[14]   Fritz BA, Chen Y, Murray-Torres TM, et al. Using machine learning techniques to develop forecasting algorithms for postoperative complications: protocol for a retrospective study. *BMJ Open* 2018; 8: e020124.

[15]   Westerman K, Reaver A, Roy C, et al. Longitudinal analysis of biomarker data from a personalized nutrition platform in healthy subjects. *Sci Rep*; 8. Published Online First: December 2018. doi: 10.1038/s41598-018-33008-7.

[16]   Minich DM, Bland JS. Personalized lifestyle medicine: relevance for nutrition and lifestyle recommendations. *The Scientific World Journal* 2013; 2013: 1–14. doi: 10.1155/2013/129841

[17]  Darcy AM, Louie AK, Roberts LW. Machine learning and the profession of medicine. *JAMA* 2016; 315: 551–552.

[18]   Sun D, Liu J, Xiao L, et al. Recent development of risk-prediction models for incident hypertension: An updated systematic review. *PLOS ONE* 2017; 12: e0187240.

[19]   Thornton RLJ, Glover CM, Cené CW, et al. Evaluating strategies for reducing health disparities by addressing the social determinants of health. *Health Aff* 2016; 35: 1416–1423. doi: 10.1377/hlthaff.2015.1357

[20]   Egger G, Dixon J. Beyond obesity and lifestyle: A review of 21st century chronic disease determinants. *BioMed Res Int* 2014; 2014: 1–12. doi: 10.1155/2014/731685

[21]   Gastil R. The determinants of human behavior. *Am Anthropol* 1961; 63: 1281–1291. http://www.jstor.org/stable/666861 (Accessed 5 Mar 2019).

[22]   Szyf M, McGowan P, Meaney MJ. The social environment and the epigenome. *Environ Mol Mutagen* 2008; 49: 46–60. doi: 10.1002/em.20357

[23] Dagum P. Digital biomarkers of cognitive function. *NPJ Digit Med*; 1. Published Online First: December 2018. doi: 10.1038/s41746-018-0018-4.

[24] Shin EK, Mahajan R, Akbilgic O, et al. Sociomarkers and biomarkers: predictive modeling in identifying pediatric asthma patients at risk of hospital revisits. *NPJ Digit Med*; 2018; 1: 50. doi:10.1038/s41746-018-0056-y

[25] Billings J, Blunt I, Steventon A, et al. Development of a predictive model to identify inpatients at risk of re-admission within 30 days of discharge (PARR-30). *BMJ Open* 2012; 2: e001667. doi: 10.1136/bmjopen-2012-001667

[26] Whelton PK, Carey RM, Aronow WS, et al. 2017 ACC/AHA/AAPA/ABC/ACPM/AGS/APhA/ASH/ASPC/NMA/PCNA Guideline for the prevention, detection, evaluation, and management of High Blood Pressure in Adults. *J Am Coll Cardiol* 2018; 71: e127–e248. doi: 10.1161/HYP.0000000000000065

[27] Williams B, Mancia G, Spiering W, et al. 2018 ESC/ESH Guidelines for the management of arterial hypertension. *Eur Heart J* 2018; 39: 3021–3104. doi: 10.1093/eurheartj/ehy339

[28] Quorum Review IRB: independent ethics review board. *Quorum Review IRB*, https://www.quorumreview.com/ (accessed 6 December 2017).

[29] Kuhn M, Johnson K. *Applied Predictive Modeling*. 5th ed. Springer, 2016.

[30] Ettehad D, Emdin CA, Kiran A, et al. Blood pressure lowering for prevention of cardiovascular disease and death: a systematic review and meta-analysis. *Lancet* 2016; 387: 957–967. doi: 10.1016/S0140-6736(15)01225-8

[31] Thomopoulos C, Parati G, Zanchetti A. Effects of blood pressure lowering on outcome incidence in hypertension. 1. Overview, meta-analyses, and meta-regression analyses of randomized trials. *J Hypertens (Los Angel)* 2014; 32: 2285–2295. doi:10.1097/HJH.0000000000000447

[32] Milman T, Joundi RA, Alotaibi NM, et al. Clinical inertia in the pharmacological management of hypertension. *Medicine (Baltimore)*; 97. Published Online First: 22 June 2018. doi: 10.1097/MD.0000000000011121.

[33] Ogedegbe G. Barriers to optimal hypertension control. *J of Clin Hypertens (Greenwich)* 2008; 10: 644–646. doi:10.1111/j.1751-7176.2008.08329.x

[34] Breiman L. Random forests. *Machine Learning* 2001; 45: 5–32. https://doi.org/10.1023/A:1010933404324

[35] scikit-learn developers. 3.2. Tuning the hyper-parameters of an estimator. *scikit-learn*, https://scikit-learn.org/stable/modules/grid_search.html (2007, accessed 31 May 2019).

[36] Liu M-X, Chen X, Chen G, et al. A Computational framework to infer human disease-associated long noncoding RNAs. *PLoS ONE* 2014; 9: e84408. doi: 10.1371/journal.pone.0084408

[37] Doytchinova IA, Flower DR. VaxiJen: a server for prediction of protective antigens, tumour antigens and subunit vaccines. *BMC Bioinformatics* 2007; 7. doi: 10.1186/1471-2105-8-4

[38] Moore RG, Brown AK, Miller MC, et al. The use of multiple novel tumor biomarkers for the detection of ovarian carcinoma in patients with a pelvic mass. *Gynecol Oncol* 2008; 108: 402–408. doi: 10.1016/j.ygyno.2007.10.017

[39] Airola A, Pahikkala T, Waegeman W, et al. An experimental comparison of cross-validation techniques for estimating the area under the ROC curve. *Comput Stat & Data Anal* 2011; 55: 1828–1844. doi: 10.1016/j.csda.2010.11.018

[40]   Lundberg SM, Lee SI. Consistent feature attribution for tree ensembles. https://arxiv.org/abs/1706.06060 (2017, accessed 19 November 2018).

[41]   Shaban-Nejad A, Michalowski M, Buckeridge DL. Health intelligence: how artificial intelligence transforms population and personalized health. *NPJ Digit Med*; 1. Published Online First: December 2018. doi: 10.1038/s41746-018-0058-9.

[42]   Alba AC, Agoritsas T, Walsh M, et al. Discrimination and Calibration of Clinical Prediction Models: Users' Guides to the Medical Literature. *JAMA* 2017; 318: 1377. doi: 10.1001/jama.2017.12126

[43]   Moons KGM, Kengne AP, Woodward M, et al. Risk prediction models: I. Development, internal validation, and assessing the incremental value of a new (bio)marker. *Heart* 2012; 98: 683–690. doi: 10.1136/heartjnl-2011-301246

[44]   Moons KGM, Kengne AP, Grobbee DE, et al. Risk prediction models: II. External validation, model updating, and impact assessment. *Heart* 2012; 98: 691–698. doi: 10.1136/heartjnl-2011-301247

**FIGURE LEGENDS**

**Figure 1**: Receiver Operator Characteristics (ROC) curves for machine learning model predicting systolic change (SC) and a model predicting systolic change without use of ongoing blood pressure data (SC-APP).

**Figure 2**: Shapley values illustrate how explanatory variables contribute to success meeting the response variable (improvement in systolic blood pressure ≥ 10 mmHg). The feature list down the y-axis is in order of contribution to the model (most to least). Each dot represents the value for one participant. SBP change and DBP change are the difference in measurements from baseline to the end of the 28-day training period.

**Figure 3**: SHAP values for explanatory variables for 2 participants. The SHAP value plotted on the y-axis indicates that amount the variable positively or negatively contributes to the prediction of success (the output value). The probability threshold (output value that assigns a prediction of success) is 0.66.
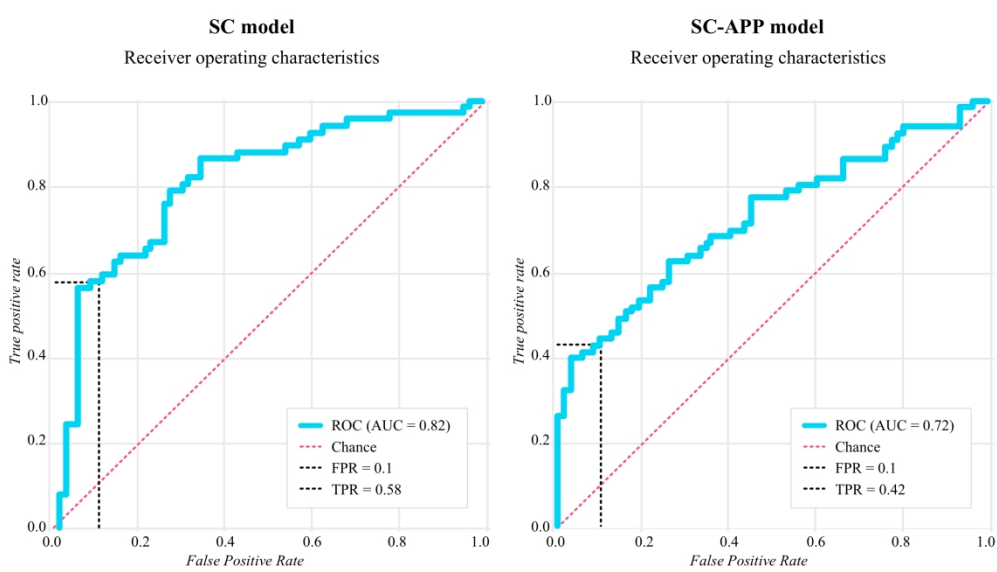
Figure 1: Receiver Operator Characteristics (ROC) curves for machine learning model predicting systolic change (SC) and a model predicting systolic change without use of ongoing blood pressure data (SC-APP).
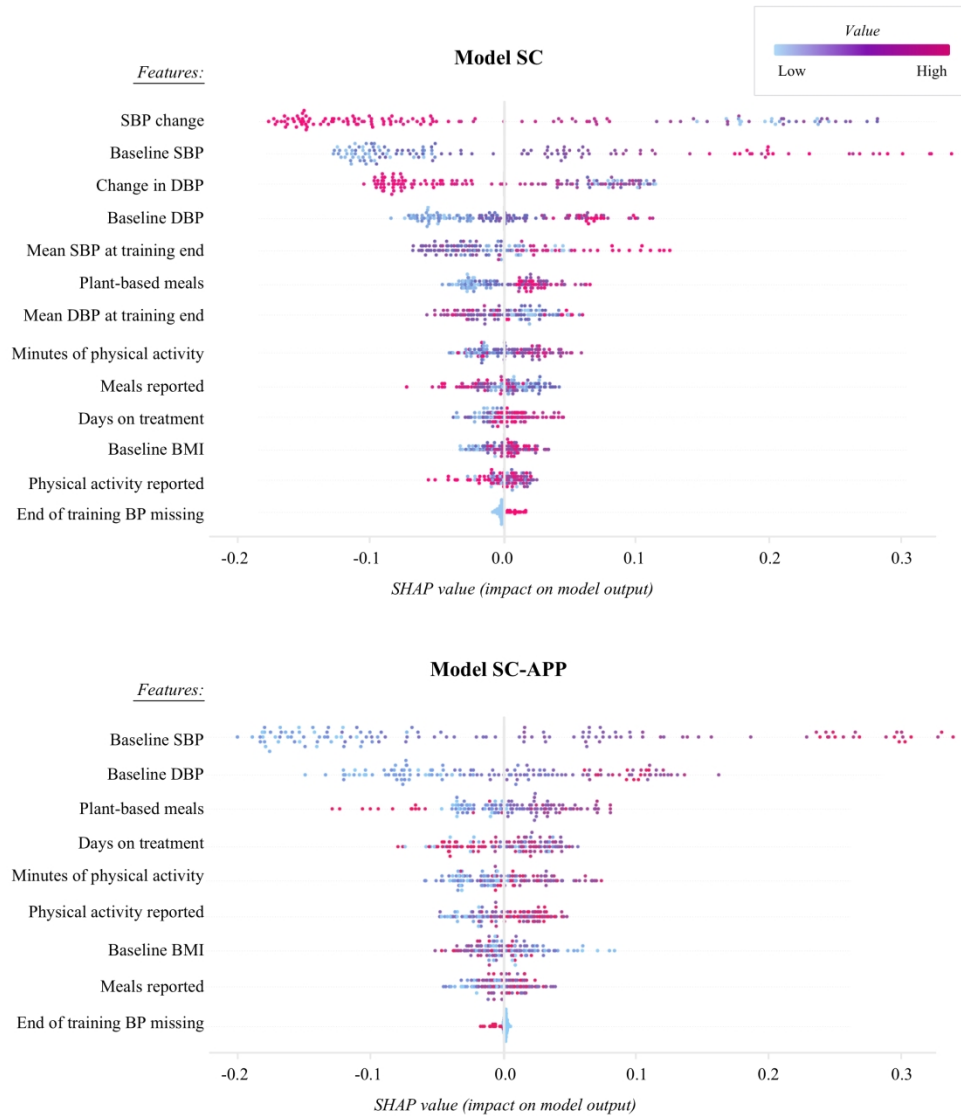
934x608mm (72 x 72 DPI)

Figure 2: Shapley values illustrate how explanatory variables contribute to success meeting the response variable (improvement in systolic blood pressure ≥ 10 mmHg). The feature list down the y-axis is in order of contribution to the model (most to least). Each dot represents the value for one participant. SBP change and DBP change are the difference in measurements from baseline to the end of the 28-day training period.

922x1058mm (72 x 72 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
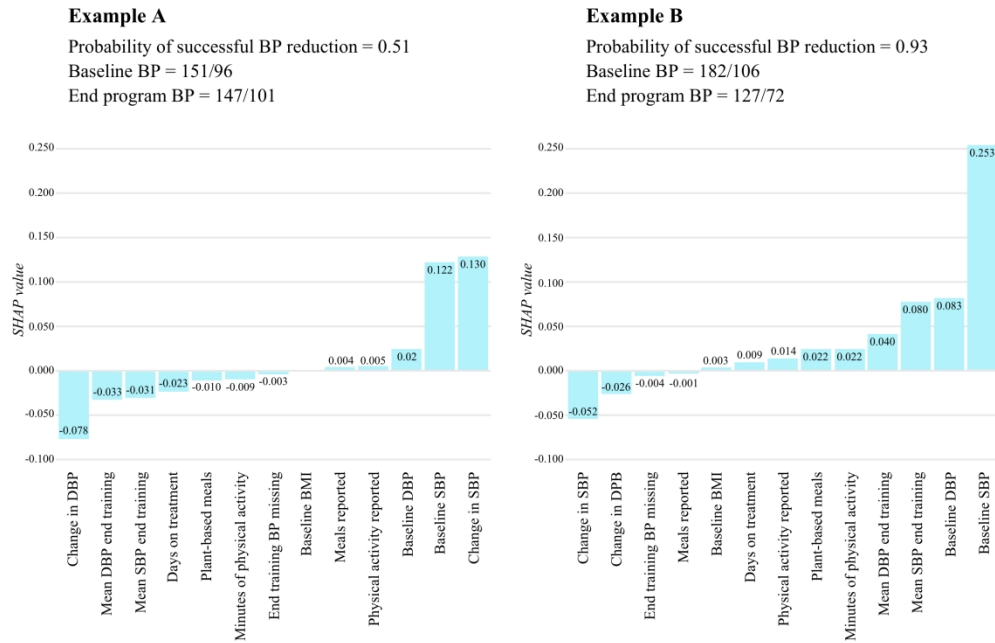21
22
23
24
25
26
27
28
29
30



Figure 3: SHAP values for explanatory variables for 2 participants. The SHAP value plotted on the y-axis indicates that amount the variable positively or negatively contributes to the prediction of success (the output value). The probability threshold (output value that assigns a prediction of success) is 0.66.

934x680mm (72 x 72 DPI)

31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

# Reporting checklist for prediction model development and validation study.

Based on the TRIPOD guidelines.

## Instructions to authors

Complete this checklist by entering the page numbers from your manuscript where readers will find each of the items listed below.

Your article may not currently address all the items on the checklist. Please modify your text to include the missing information. If you are certain that an item does not apply, please write "n/a" and provide a short explanation.

Upload your completed checklist as an extra file when you submit to a journal.

In your methods section, say that you used the TRIPOD reporting guidelines, and cite them as:

Collins GS, Reitsma JB, Altman DG, Moons KG. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): The TRIPOD statement.

|  | | Reporting Item | Page Number |
|---|---|---|---|
| | #1 | Identify the study as developing and / or validating a multivariable prediction model, the target population, and the outcome to be predicted. | 2 |
| | #2 | Provide a summary of objectives, study design, setting, participants, sample size, predictors, outcome, statistical analysis, results, and conclusions. | 2 |
| | #3a | Explain the medical context (including whether diagnostic or prognostic) and rationale for developing or validating the multivariable prediction model, including references to existing models. | 3-4 |
| | #3b | Specify the objectives, including whether the study describes the development or validation of the model or both. | 4 |
| Source of data | #4a | Describe the study design or source of data (e.g., randomized trial, cohort, or registry data), separately for the development and validation data sets, if applicable. | 4-5 |

| | | | |
|---|---|---|---|
| | #4b | Specify the key study dates, including start of accrual; end of accrual; and, if applicable, end of follow-up. | n/a |
| Participants | #5a | Specify key elements of the study setting (e.g., primary care, secondary care, general population) including number and location of centres. | 4 |
| | #5b | Describe eligibility criteria for participants. | 5 |
| | #5c | Give details of treatments received, if relevant | 4 |
| Outcome | #6a | Clearly define the outcome that is predicted by the prediction model, including how and when assessed. | 5-6 |
| | #6b | Report any actions to blind assessment of the outcome to be predicted. | n/a |
| Predictors | #7a | Clearly define all predictors used in developing or validating the multivariable prediction model, including how and when they were measured | 5-7 |
| | #7b | Report any actions to blind assessment of predictors for the outcome and other predictors. | n/a |
| Sample size | #8 | Explain how the study size was arrived at. | 5 |
| Missing data | #9 | Describe how missing data were handled (e.g., complete-case analysis, single imputation, multiple imputation) with details of any imputation method. | n/a |
| Statistical analysis methods | #10a | If you are developing a prediction model describe how predictors were handled in the analyses. | 5-7 |
| | #10b | If you are developing a prediction model, specify type of model, all model-building procedures (including any predictor selection), and method for internal validation. | 6-8 |
| | #10c | If you are validating a prediction model, describe how the predictions were calculated. | 7-8 |
| | #10d | Specify all measures used to assess model performance and, if relevant, to compare multiple models. | 7-8 |
| | #10e | If you are validating a prediction model, describe any model updating (e.g., recalibration) arising from the validation, if done | 7-8 |

| | | | |
|---|---|---|---|
| Risk groups | #11 | Provide details on how risk groups were created, if done. | n/a |
| Development vs. validation | #12 | For validation, identify any differences from the development data in setting, eligibility criteria, outcome, and predictors. | 8 |
| Participants | #13a | Describe the flow of participants through the study, including the number of participants with and without the outcome and, if applicable, a summary of the follow-up time. A diagram may be helpful. | 9 |
| | #13b | Describe the characteristics of the participants (basic demographics, clinical features, available predictors), including the number of participants with missing data for predictors and outcome. | 9-10 |
| | #13c | For validation, show a comparison with the development data of the distribution of important variables (demographics, predictors and outcome). | 9-10 |
| Model development | #14a | If developing a model, specify the number of participants and outcome events in each analysis. | 9-10 |
| | #14b | If developing a model, report the unadjusted association, if calculated between each candidate predictor and outcome. | 10-11 |
| Model specification | #15a | If developing a model, present the full prediction model to allow predictions for individuals (i.e., all regression coefficients, and model intercept or baseline survival at a given time point). | 10-11 |
| | #15b | If developing a prediction model, explain how to the use it. | 10-11 |
| Model performance | #16 | Report performance measures (with CIs) for the prediction model. | 10-11 |
| Model-updating | #17 | If validating a model, report the results from any model updating, if done (i.e., model specification, model performance). | n/a |
| Limitations | #18 | Discuss any limitations of the study (such as nonrepresentative sample, few events per predictor, missing data). | 12-13 |
| Interpretation | #19a | For validation, discuss the results with reference to performance in the development data, and any other validation data | 11-12 |
| | #19b | Give an overall interpretation of the results, considering objectives, limitations, results from similar studies, and other relevant evidence. | 14 |

| | | | | |
|---|---|---|---|---|
| Implications | #20 | Discuss the potential clinical use of the model and implications for future research | | 13-14 |
| Supplementary information | #21 | Provide information about the availability of supplementary resources, such as study protocol, Web calculator, and data sets. | | 14-15 |
| Funding | #22 | Give the source of funding and the role of the funders for the present study. | | 14 |

The TRIPOD checklist is distributed under the terms of the Creative Commons Attribution License CC-BY. This checklist was completed on 27. March 2019 using https://www.goodreports.org/, a tool made by the EQUATOR Network in collaboration with Penelope.ai