

PERSONALIZED HYPOTHESIS TESTS FOR DETECTING MEDICATION RESPONSE IN PARKINSON DISEASE PATIENTS USING iPhone SENSOR DATA

ELIAS CHAIBUB NETO*, BRIAN M. BOT, THANNEER PERUMAL, LARSSON OMBERG, JUSTIN GUINNEY, MIKE KELLEN, ARNO KLEIN, STEPHEN H. FRIEND, ANDREW D. TRISTER

Sage Bionetworks, 1100 Fairview Avenue North, Seattle, Washington 98109, USA

** corresponding author e-mail: elias.chaibub.neto@sagebase.org*

We propose hypothesis tests for detecting dopaminergic medication response in Parkinson disease patients, using longitudinal sensor data collected by smartphones. The processed data is composed of multiple features extracted from active tapping tasks performed by the participant on a daily basis, before and after medication, over several months. Each extracted feature corresponds to a time series of measurements annotated according to whether the measurement was taken before or after the patient has taken his/her medication. Even though the data is longitudinal in nature, we show that simple hypothesis tests for detecting medication response, which ignore the serial correlation structure of the data, are still statistically valid, showing type I error rates at the nominal level. We propose two distinct personalized testing approaches. In the first, we combine multiple feature-specific tests into a single union-intersection test. In the second, we construct personalized classifiers of the before/after medication labels using all the extracted features of a given participant, and test the null hypothesis that the area under the receiver operating characteristic curve of the classifier is equal to $1/2$. We compare the statistical power of the personalized classifier tests and personalized union-intersection tests in a simulation study, and illustrate the performance of the proposed tests using data from mPower Parkinsons disease study, recently launched as part of Apples ResearchKit mobile platform. Our results suggest that the personalized tests, which ignore the longitudinal aspect of the data, can perform well in real data analyses, suggesting they might be used as a sound baseline approach, to which more sophisticated methods can be compared to.

Keywords: personalized medicine, hypothesis tests, sensor data, remote monitoring, Parkinson

1. Introduction

Parkinson disease is a severe neurodegenerative disorder of the central nervous system caused by the death of dopamine-generating cells in the midbrain. The disease has considerable worldwide morbidity and is associated with substantial decrease in the quality of life of the patients (and their caregivers), decreased life expectancy, and high costs related to care. Early symptoms in the motor domain include shaking, rigidity, slowness of movement and difficulty for walking. Later symptoms include issues with sleeping, thinking and behavioral problems, depression, and finally dementia in the more advanced stages of the disease. Treatments are usually based on levodopa and dopamine agonist medications. Nonetheless, as the disease progresses, these drugs often become less effective, while still causing side effects, including involuntary twisting movements (dyskinesias). Statistical approaches aiming to determine if a given patient responds to medication have key practical importance as they can help the physician in making more informed treatment recommendations for a particular patient.

In this paper we propose personalized hypothesis tests for detecting medication response in Parkinson patients, using longitudinal sensor data collected by iPhones. Remote monitoring of Parkinson patients, based on active tasks delivered by smartphone applications, is an active research field.¹ Here we illustrate the application of our personalized tests using sensor data

collected by the mPower study, recently launched as part of Apple's ResearchKit^{2,3} mobile platform. The active tests implemented in the mPower app include tapping, voice, memory, posture and gait tests, although in this paper we focus on the tapping data only. During a tapping test the patient is asked to tap two buttons on the iPhone screen alternating between two fingers on the same hand for 20 seconds. Raw sensor data collected during a single test is given by a time series of the screen x-y coordinates on each tap. Processed data corresponds to multiple features extracted from the tapping task, such as the number of taps and the mean inter-tapping interval. Since the active tests are performed by the patient on a daily basis, before and after medication, over several months, the processed data corresponds to time series of feature measurements annotated according to whether the measurement was taken before or after the patient has taken his/her medication. Though others have investigated the feasibility of monitoring medication response in Parkinson patients using smartphone sensor data, this previous work did not focus on the individual effects that medications have, but rather focused on the classification on a population level.⁴

The first step in analyzing these data is to show that simple feature-specific tests, which ignore the serial correlation in the extracted features, are statistically valid (the distribution of the p-values for tests applied to data generated under the null hypothesis is uniform). This condition guarantees that the tests are exact, that is, the type I error rates match the nominal levels, so that our inferences are neither conservative nor liberal. In other words, if we adopt a significance level cutoff of α , the probability that our tests will incorrectly reject the null when it is actually true is given by α .

Even though the simple feature-specific tests are valid procedures for testing for medication response, in practice, we have multiple features and need to combine them into a single decision procedure. The second main contribution of this paper is to propose two distinct approaches to combine all the extracted features into a single hypothesis test. In the first, and most standard approach, we combine simple tests, applied to each one of our extracted features, into a single union-intersection test. Although simple to implement, scalable, and computationally efficient, this approach requires multiple testing correction, which might become burdensome when the number of extracted features is large. In order to circumvent this potential issue, our second approach is to construct personalized classifiers of the before/after medication labels using all the extracted features of a given patient, and test the null hypothesis that the area under the receiver operating characteristic curve (AUROC) of the classifier is equal to $1/2$ (in which case the patient's extracted features are unable to predict the before/after medication labels, implying that the patient does not respond to the medication). A slight disadvantage of the classifier approach, compared to the union-intersection tests, is the larger computational cost (especially for classifiers that require tuning parameter optimization by cross-validation) involved in the classifier training. In any case, the increased computational demand is by no means a limiting factor for the application of the approach.

The rest of this paper is organized as follows. In Section 2 we present our personalized tests, discuss their statistical validity, and perform a power study comparison. In Section 3 we illustrate the application of our tests to the tapping data of the mPower study. Finally, in Section 4 we discuss our results.

2. Methods

2.1. *Notation and a few preliminary comments on the data*

Throughout this paper, we let x_{kt} , $k = 1, \dots, p$, $t = 1, \dots, n$, represent the measurement of feature k at time point t , and let $y_t = \{b, a\}$, represent the binary outcome variable, corresponding to the before/after medication label, where b and a stand for “before” and “after” medication, respectively.

Even though the participants were asked to perform the active tasks 3 times per day, one before the medication, one after, and one at any other time of their choice, participants did not always follow the instructions correctly. As a result, the data is non-standard, with variable number of daily tasks (sometimes fewer, sometimes greater than 3 tasks per day), and variable timing relative to medication patterns (e.g., bbabba... , aaabbb... , instead of bababa...). Furthermore, the data also contains missing medication labels, as sometimes, a participant performed the active task but did not report whether the task was taken before or after medication. In our analysis we restrict our attention to data collected before and after medication only. Hence, for each participant, the number of data points used in our tests is given by $n = n_b + n_a$, where n_b and n_a correspond, respectively, to the number of before/after medication labels.

2.2. *On the statistical validity of personalized tests which ignore the autocorrelation structure of the data*

It is common knowledge that the t-test, the Wilcoxon rank-sum test, and other two-sample problem tests, suffer from inflated type I error rates in the presence of dependency. We point out, however, that this can happen when the data within each group is dependent, but the two groups are themselves statistically independent. When the data from both groups is sampled jointly from the same multivariate distribution, the dependency of the data might no longer be an issue. Figure 1 provides an illustrative example with t-tests applied to simulated data.

The t-test’s assumption of independence (within and between the groups’ data) is required in order to make the analytical derivation of the null distribution feasible. It doesn’t mean the test will always generate inflated type I error rates in the presence of dependency (as illustrated in Figure 1f). As a matter of fact, a permutation test based on the t-test statistic is valid if the group labels are exchangeable under the null,⁵ even when the data is statistically dependent. Exchangeability⁶ captures a notion of symmetry/similarity in the data, without requiring independence. On the examples presented in Figure 1, the group labels are exchangeable on panels a and c as illustrated by the symmetry/similarity of the data between groups 1 and 2 at each row of the heatmaps. For panel b, on the other hand, the lack of symmetry between the groups on each row illustrates that the group labels are not exchangeable.

In the context of our personalized tests, the before/after medication labels are exchangeable under the null of no medication response, even though the measurements of any extracted feature are usually serially correlated. Note that the exchangeability is required for the medication labels, and not for the feature measurements, which are not exchangeable due to their serial correlation. Figure 2 illustrates this point, showing the symmetry/similarity of the separate time series for the before and after medication data.

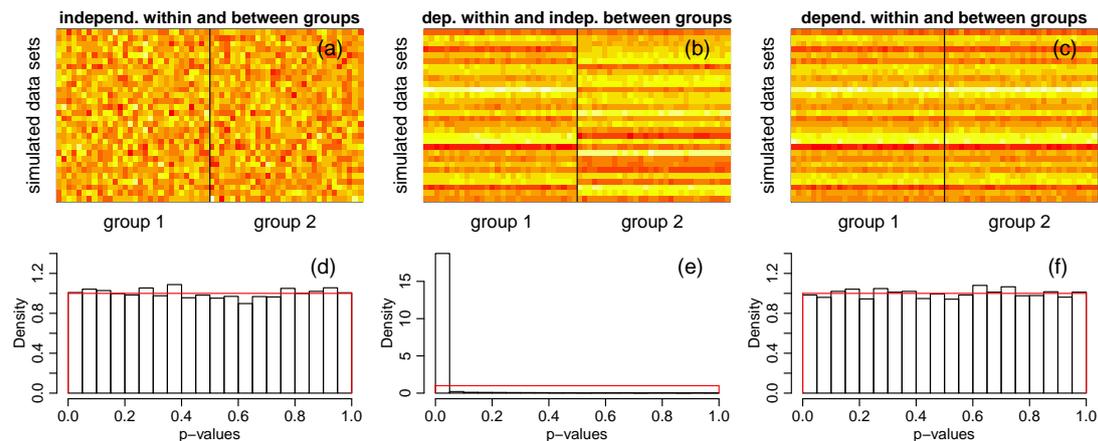


Fig. 1. The effect of data dependency on the t-test. Panels a, b, and c show heatmaps of the simulated data. Columns are split between group 1 and 2, and each row corresponds to one simulated null data set (we show the top 30 simulations only). Bottom panels show the p-value distributions for 10,000 tests applied to null data simulated according to: (i) $\mathbf{x}_1 \sim N_{30}(\boldsymbol{\mu}_1, \mathbf{I})$ and $\mathbf{x}_2 \sim N_{30}(\boldsymbol{\mu}_2, \mathbf{I})$ with $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = \mathbf{0}$ (panel d); (ii) $\mathbf{x}_1 \sim N_{30}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma})$ and $\mathbf{x}_2 \sim N_{30}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma})$, where $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = \mathbf{0}$ and $\boldsymbol{\Sigma}$ is a correlation matrix with off-diagonal elements equal to $\rho = 0.95$ (panel e); and (iii) $(\mathbf{x}_1, \mathbf{x}_2)^t \sim N_{60}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, with $\boldsymbol{\mu} = (\boldsymbol{\mu}_1, \boldsymbol{\mu}_2)^t = \mathbf{0}$ and $\boldsymbol{\Sigma}$ as before (panel f). The density of the Uniform[0, 1] distribution is shown in red. Panel d shows that under the standard assumptions of the t-test, the p-value distribution under the null is (as expected) uniform. Panel e shows the p-value distribution for strongly dependent data, showing highly inflated type I error rates, even though the data was simulated according to t-test's null hypothesis that $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$. Panel b clarifies why that is the case. For each row (i.e., simulated data set), the data tends to be quite homogeneous inside each group, but quite distinct between the groups. Because on each simulation we sample the data vectors \mathbf{x}_1 and \mathbf{x}_2 from a multivariate normal distribution with a very strong correlation structure, all elements in the \mathbf{x}_1 vector tend to be close to each other, and all elements in \mathbf{x}_2 tend to be similar to each other. However, because \mathbf{x}_1 and \mathbf{x}_2 are sampled independently from each other, their values tend to be distinct. In combination, the small variability in each group vector together with the difference in their means leads to high test statistic values and small p-values. Panel f shows the p-value distribution for strongly dependent data, when sampled jointly. In this case, the distribution is uniform. Panel c clarifies why. Now, each row tends to be entirely homogeneous (within and between groups), since the joint sampling of \mathbf{x}_1 and \mathbf{x}_2 makes all elements in both vectors quite similar to each other, so that the difference in their means tends to be small.

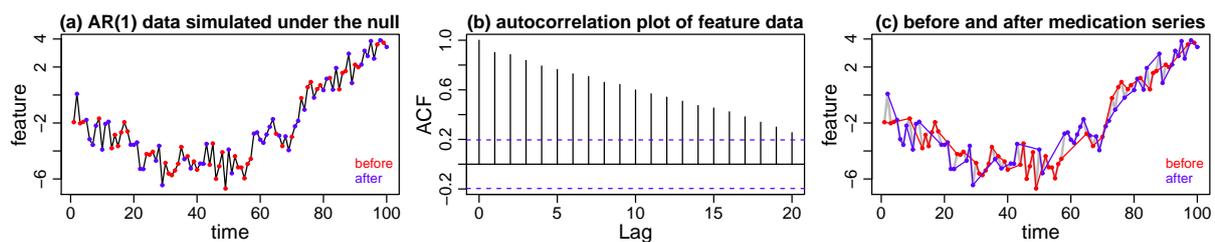


Fig. 2. Exchangeability of the before/after medication labels in time series data. Panel a shows a feature, simulated from an AR(1) process, under the null hypothesis that the patient does not respond to medication. In this case, the before medication (red dots) and after medication (blue dots) labels are randomly assigned to the feature measurements. Panel b shows an autocorrelation plot of the feature data. Panel c shows the separate “before medication” (red) and “after medication” (blue) series. Note the symmetry/similarity of the two series. Clearly, under the null hypothesis that the patient does not respond to medication, the medication labels are exchangeable, since shuffling of the before/after medication labels would not destroy the serial correlation structure or the trend of the series.

Hence, even though our longitudinal data violates the independence assumption of the t-test, the permutation test based on the t-test statistic is still valid. Of course the same argument is valid for permutation tests based on other test statistics. Figure 3 illustrates this point, with permutation tests based on the t-test and Wilcoxon rank-sum test. Panel a shows the original data. Red and blue dots represent measurements before and after medication, respectively. The grey dots represent data collected at another time or where the medication label is missing. Panel b shows one realization of a random permutation of the before/after medication labels. In order to generate a permutation null distribution, we perform a large number of random label shuffles, and for each one, we evaluate the adopted test statistic in the permuted data. Panels c and d show the permutation null distributions generated from 10,000 random permutations of the medication labels based, respectively, on the t-test and on the Wilcoxon rank-sum test statistics. The red curve on panel c shows the analytical density of a

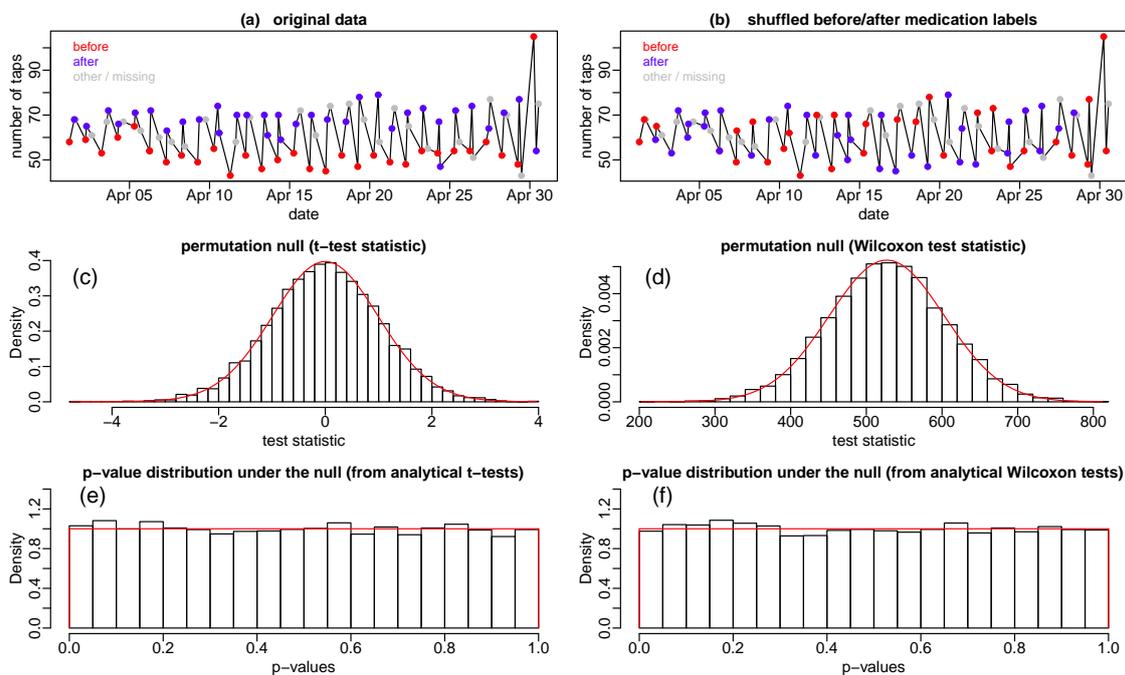


Fig. 3. Personalized tests for the null hypothesis that the patient does not respond to medication, according to a single feature (number of taps), and based on t-test and Wilcoxon rank-sum test statistics. Panel a shows the data on the number of taps for one of mPower’s study participants during April 2015. Panel b shows one realization of a random permutation of the before/after medication labels. Panel c shows the permutation null distribution based on the t-test statistic. The red curve represents the analytical density of a t-test, namely, a t-distribution with $n_b + n_a - 2$ degrees of freedom. Panel d shows the permutation null distribution based on Wilcoxon rank-sum test statistic. The red curve shows the density of the normal asymptotic approximation for Wilcoxon’s test, namely, a Gaussian density with mean, $n_b n_a / 2$, and variance, $n_b n_a (n_b + n_a + 1) / 12$. In this example $n_b = 31$ and $n_a = 34$. Panels e and f show the analytical p-value distributions under the null. Results are based on 10,000 null data sets. Each null data set was generated as follows: (i) randomly sample the number of “before” medication labels, n_b , from the set $\{10, 11, \dots, n - 10\}$, where n is the total number of measurements; (ii) compute the number of “after” medication labels as $n_a = n - n_b$; and (iii) randomly assign the “before medication” and “after medication” labels to the number of taps measurements. The plots show the histograms of the p-values derived from the application of t-tests (panel e) and Wilcoxon’s tests (panel f) to each of the null data sets. The density of the Uniform $[0, 1]$ distribution is shown in red.

t-test for the data on panel a, while the red curve on panel d shows the density of the normal asymptotic approximation for the Wilcoxon test represented in panel b (we show the normal approximation density, since the exact null distribution is discrete). The close similarity of the permutation and analytical distributions suggests that, in practice, we can use the analytical p-values of t-tests or Wilcoxon rank-sum tests instead of the permutation p-values. Panels e and f further corroborate this point, showing uniform distributions for the analytical p-values of t-tests and Wilcoxon tests respectively, derived from 10,000 distinct null data sets.

2.3. Personalized union-intersection tests

In order to combine the feature-specific tests, H_{0k} : the patient does not respond to the medication, according to feature k , versus H_{1k} : the patient responds to the medication, across all extracted features, we construct the union-intersection test,

$$H_0 : \cap_{k=1}^p H_{0k} \quad \text{versus} \quad H_1 : \cup_{k=1}^p H_{1k} , \quad (1)$$

where, in words, we test the null hypothesis that the patient does not respond to medication, for all p features, versus the alternative hypothesis that he/she responds to medication according to at least one of the features. Under this test, we reject H_0 if the p-value of at least one of the feature-specific tests is small. Hence, the p-value for the union-intersection test corresponds to the smallest p-value (across all p tests) after multiple testing correction.

We implement union-intersection tests based on the t-test and Wilcoxon rank-sum test statistics. We adopt the Benjamini-Hochberg approach⁷ for multiple testing correction. As detailed in Figure 4, the union-intersection test tends to be slightly conservative when applied to correlated features.

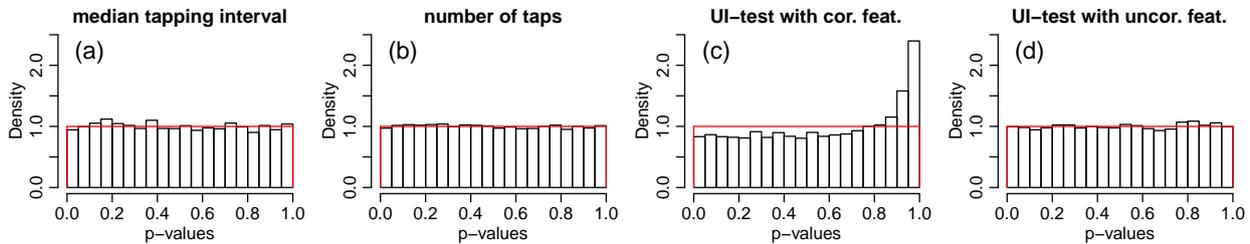


Fig. 4. P-value distributions for the union-intersection test under the null. Results are based on 10,000 null data sets generated by randomly permuting the before/after medication labels. Panels a and b show the p-value distributions from, H_{0k} , for 2 out of the 24 features combined in the union-intersection test. We observed uniform distributions for all features, but report just 2 here due to space constraints. Panel c shows the p-value distribution of the union-intersection test using Benjamini-Hochberg correction. The skewness of the distribution towards larger p-values indicates that the test is conservative, meaning that at a nominal significance level α the probability of rejecting the null when it is actually true is smaller than α . Since the p-value distributions of the features are uniform (panels a and b), the skewness of the union-intersection p-value distribution is clearly due to the multiple testing correction. We experimented with other procedures, but the Benjamini-Hochberg correction was the least conservative one. Panel d reports the p-value distribution for the union-intersection test using a shuffled version of the feature data. Note that by shuffling the data, we destroy the correlation among the features, so that the now uniform distribution suggests that the violation of the independence assumption (required by the Benjamini-Hochberg approach) seems to be the reason for the slightly conservative behavior of the union-intersection test. Results were generated using Wilcoxon's test.

2.4. Personalized classifier tests

An alternative approach to combine the information across multiple features into a single decision procedure is to test whether a classifier trained using all the features is able to predict the before/after medication labels better than a random guess. Adopting the AUROC as a classification performance metric, we have that a prediction equivalent to a random guess would have an AUROC equal to 0.5, whereas a perfect prediction would lead to an AUROC equal to 1. Furthermore, if a classifier generates an AUROC smaller than 0.5, we only need to switch the labels in order to make the AUROC larger than 0.5. Therefore, we can test if a classifier's prediction is better than random using the one-sided test,

$$H_0 : \text{AUROC} = 1/2 \quad \text{versus} \quad H_1 : \text{AUROC} > 1/2 . \quad (2)$$

It has been shown⁸ that, when there are no ties in the predicted class probabilities used for the computation of the AUROC, the test statistic of the Wilcoxon rank-sum test (also known as the Mann-Whitney U test), U , is related to the AUROC statistic by $U = n_b n_a (1 - \text{AUROC})$ (see section 2 of reference⁹ for details). Hence, under the assumption of independence (required by Wilcoxon's test) the analytical p-value for the hypothesis test in (2) is given by the left tail probability, $P(U \leq n_b n_a (1 - \text{AUROC}))$, of Wilcoxon's null distribution. In the presence of ties, the p-value is given by the left tail of the asymptotic approximate null,

$$U \approx N \left(\frac{n_b n_a}{2}, \frac{n_b n_a (n+1)}{12} - \frac{n_b n_a}{12 n (n-1)} \sum_{j=1}^{\tau} t_j (t_j - 1) (t_j + 1) \right), \quad (3)$$

where τ is the number of groups of ties, and t_j is the number of ties in group j .⁹ Alternatively, we can get the p-value as the right tail probability of the corresponding AUROC null,

$$\text{AUROC} \approx N \left(\frac{1}{2}, \frac{n+1}{12 n_b n_a} - \frac{1}{12 n_b n_a n (n-1)} \sum_{j=1}^{\tau} t_j (t_j - 1) (t_j + 1) \right). \quad (4)$$

As before, even though the test described above assumes independence, the exchangeability of the before/after medication labels guarantees the validity of the permutation test based on the AUROC statistic. Figure 5 illustrates this point.

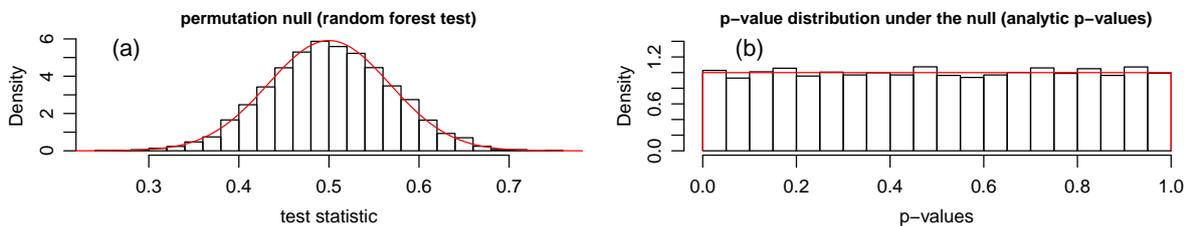


Fig. 5. Panel a shows the permutation null distribution for the personalized classifier test (based on the random forest algorithm). The red curve represents the density of the AUROC (approximate) null distribution in eq. (4). Panel b shows the p-value distribution for the (analytical) classifier test, based on 10,000 null data sets, generated by randomly permuting the before/after medication labels as described in Figure 3. The density of the Uniform[0, 1] distribution is shown in red.

2.5. Statistical power comparison

In this section we compare the statistical power of the personalized classifier test (based on the random forest¹⁰ and extra trees¹¹ classifiers) against the personalized union-intersection tests (based on t-tests and Wilcoxon rank-sum tests). We simulated feature data, x_{kt} , according to the model,

$$x_{kt} = A \cos(\pi t) + \epsilon_{kt} , \quad (5)$$

where $\epsilon_{kt} \sim N(0, \sigma_k^2)$ represents i.i.d. error terms, and the function $A \cos(\pi t)$ describes the periodic signal, with A representing the peak amplitude of the signal. Figure 6 describes the additional steps involved in the generation of the before/after medication labels.

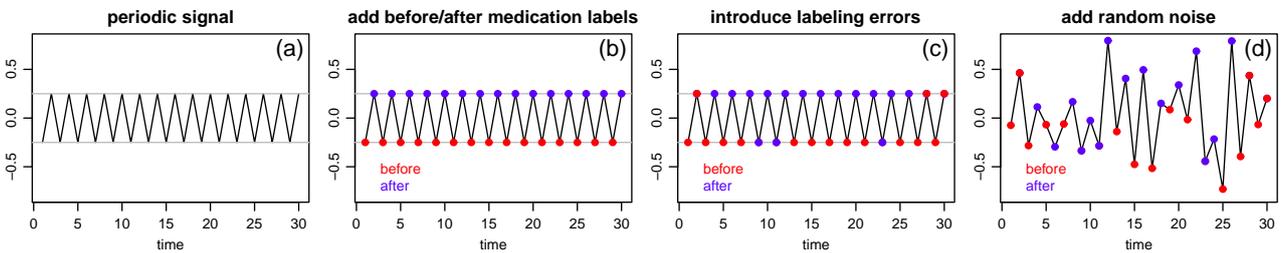


Fig. 6. Data simulation steps. First, we generate the periodic signal, $A \cos(\pi t)$, shown in panel a for $t = 1, \dots, 30$, and $A = 0.25$. Second, we assign the “before medication” label (red dots) to the negative values, and the “after medication” label (blue dots) for the positive values (panel b). Third, we introduce some labeling errors (panel c). Finally, at the fourth step (panel d), we add $\epsilon_{kt} \sim N(0, \sigma_k^2)$ error terms to the measurements.

We ran six simulation experiments, covering all combinations of sample size, $n = \{50, 150\}$, and number of features, $p = \{10, 50, 200\}$. In all simulations we adopted $A = 0.25$, mislabeling error rate of 10%, and increasing $\sigma_k = \{1.00, 1.22, 1.44, 1.67, 1.89, 2.11, 2.33, 2.56, 2.78\}$ for the first 9 features, and $\sigma_k = 3$, for $10 \leq k \leq p$. In each experiment, we generated 1,000 data sets and computed the personalized tests p-values. Figure 7 presents a comparison of the empirical power of the personalized tests as a function of the significance threshold α . For each test, the empirical power curve was estimated as the fraction of the p-values smaller than α , for a dense grid of α values varying between 0 and 1.

The results showed some clear patterns. First, the comparison of the top and bottom panels showed that for all 4 tests an increase in sample size leads to an increase in statistical power (as one would have expected). Second, the empirical power of the union-intersection tests based on Wilcoxon and t-tests was very similar in all simulation experiments, with the t-test being slightly better powered than the Wilcoxon test. This result was expected since the simulated features were generated using Gaussian errors, and the t-test is known to be slightly better powered than the Wilcoxon rank-sum test under normality. Similarly, the empirical power of the personalized classifier tests was also similar, with the extra trees algorithm tending to be slightly better powered. Third, this study shows that neither the personalized classifier nor the union-intersection approaches dominate each other. Rather, we see that for this particular data generation mechanism and simulation parameters choice, the union-intersection tests tended to be better powered than the classifier tests for smaller values of p , whereas the converse was

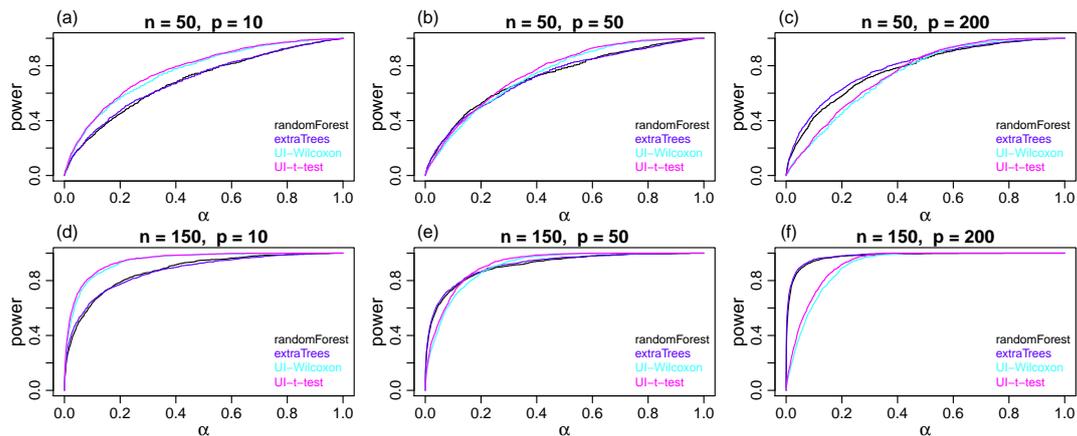


Fig. 7. Comparison of the personalized test's empirical power as a function of the significance cutoff, α .

true for larger p . Furthermore, observe that the power of the union-intersection tests tended to decrease as the number of features increased (note the slight decrease in the slope of the cyan and pink curves as we go from the panels on the left to the panels on the right). The tests based on the classifiers, on the other hand, tended to get better powered as the number of features increased (note the respective increase in the slope of the blue and black curves).

The observed decrease in power of the union-intersection tests might be explained by the increased burden caused by the multiple testing correction required by these tests. On the other hand, the personalized classifier tests are not plagued by the multiple testing issue since all features are simultaneously accounted for by the classifiers. Furthermore, in situations where none of the features is particularly informative, the classifiers might still be able to better aggregate information across the multiple noisy features.

3. Real data illustrations

In this section we illustrate the performance of our personalized hypothesis tests, based on tapping data collected by the mPower study between 03/09/2015 (date the study opened) and 06/25/2015. We restrict our analyzes to 57 patients, who performed at least 30 tapping tasks before medication, as well as 30 or more tasks after medication.

Figure 8 presents the results. Panel a shows the number of tapping tasks (before and after medication) performed by each participant. Panel b reports the AUROC scores for 4 classifiers (random forest, logistic regression with elastic-net penalty,¹² logistic regression, and extra trees). Panel c presents the p-values for the respective personalized classifier tests^a, as well as for 2 personalized union-intersection tests (based on t- and Wilcoxon rank-sum tests).

Panel b shows that the tree-based classifier tests (random forest and extra trees) showed comparable performance across all participants, whereas the regression-based approaches (elastic net and logistic regression) were sometimes comparable but sometimes strongly out-

^aThe results for the personalized classifier tests were based on 100 random splits of the data into training and test sets, using roughly half of the data for training and the other half for testing. The AUROC and p-values reported on Figure 8 b and c correspond to the median of the AUROCs and p-values across the 100 data splits.

performed by the tree-based tests. Panel c shows that the union-intersection tests, nonetheless, produced at times much smaller p-values than the classifier tests. At a significance level of 0.001 (grey horizontal line), about one quarter of the patients (leftmost quarter) respond to dopaminergic medication, according to most tests.

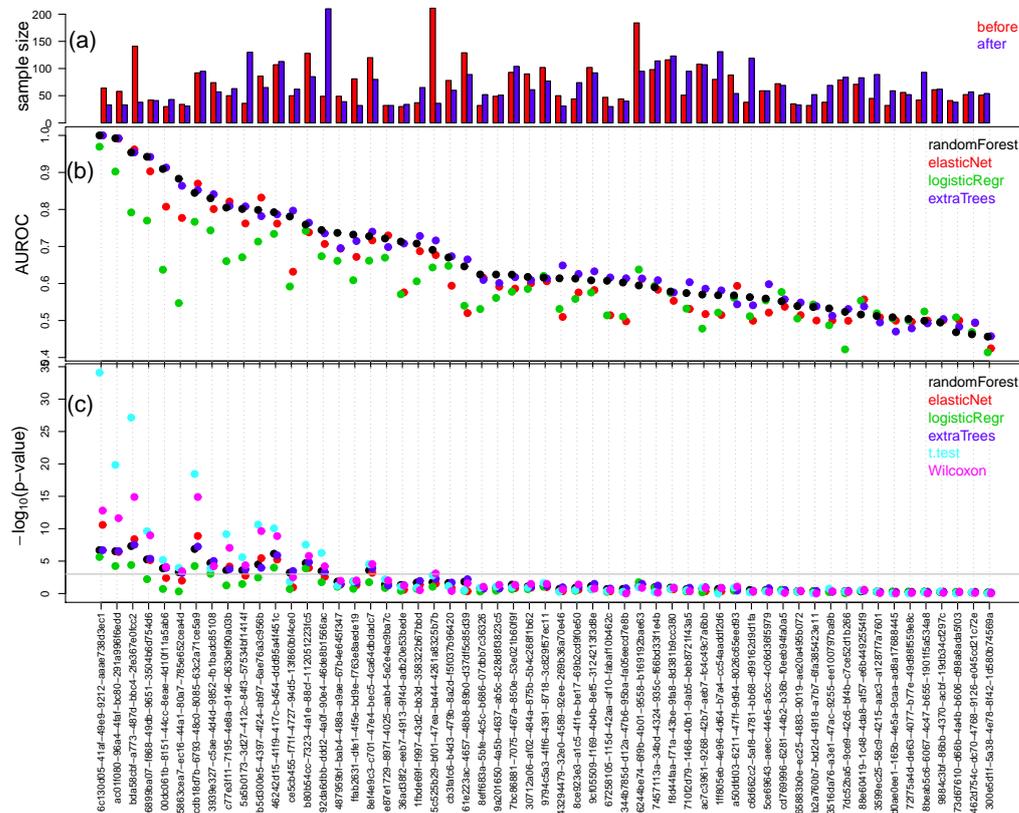


Fig. 8. Results of personalized hypothesis tests applied to the tapping data from the mPower study. Panel a shows the number of tapping tasks performed before (red) and after (blue) medication, per participant. Panel b shows the AUROC scores (across 100 random splits of the data into training and test sets) for four classifiers. Panel c shows the adjusted p-values (in negative log base 10 scale) of 4 classification tests and 2 union-intersection tests. The p-values were adjusted using Benjamini-Hochberg multiple testing correction. The grey horizontal line represents an adjusted p-value cutoff of 0.001. The participants were sorted according to the AUROC of the random forest algorithm (black dots in panel b).

In order to illustrate how our personalized tests are able to pick up meaningful signal from the extracted features, we present on Figure 9, time series plots of 2 features (number of taps and mean tapping interval) which were consistently ranked among the top 3 features for the top 3 patients on the left of Figure 8c, and compared it to the data from the bottom 3 patients on the right of Figure 8c. (For the random forest classifier test, we ranked the features according to the importance measure provided by the random forest algorithm. For the union-intersection tests, we ranked the features according to the p-values of the feature-specific tests.) Comparison of panels a to f, which show the data from patients that respond to medication (according to our tests), against panels g to l, which report the data from patients that do not respond to medication, shows that our tests can clearly pick up meaningful signal

from these two extracted features.

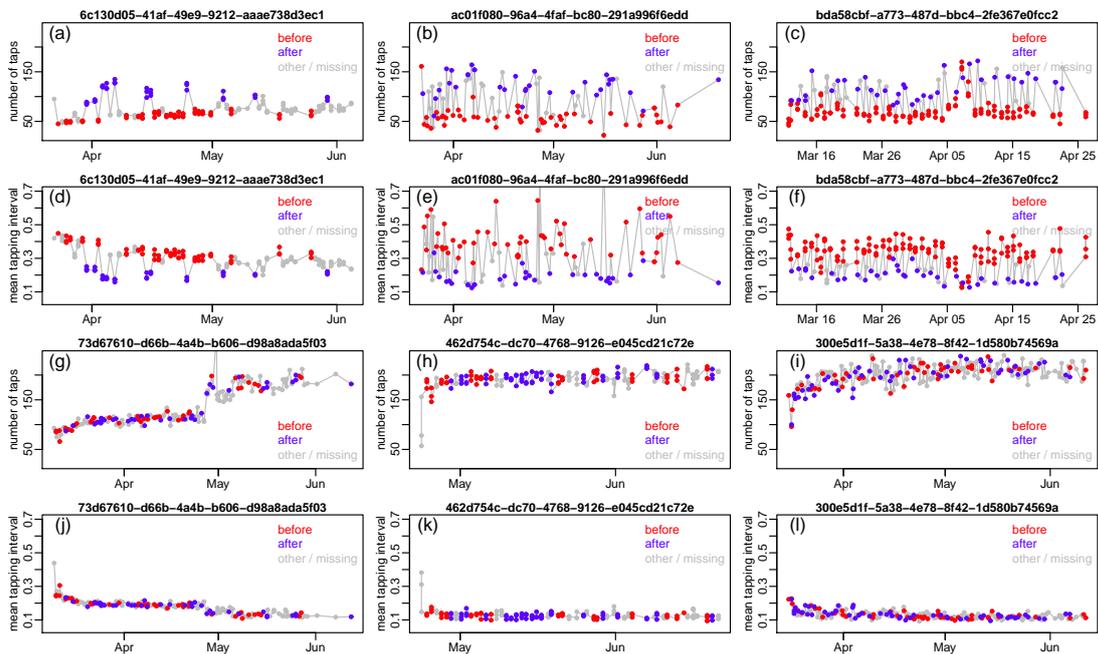


Fig. 9. Comparison of the leftmost 3 patients against the rightmost 3 patients in Figure 8c, according to 2 tapping features (number of taps and mean tapping interval). Panels a to f show the results for the top 3 patients (which respond to medication according to our tests), while panels g to l show the results for the bottom 3 patients (which don't respond to medication according to our tests).

4. Discussion

In this paper we describe personalized hypothesis tests for detecting dopaminergic medication response in Parkinson patients. We propose and compare two distinct strategies for combining the information from multiple extracted features into a single decision procedure, namely: (i) union-intersection tests; and (ii) hypothesis tests based on classifiers of the medication labels. We carefully evaluated the statistical properties of our tests, and illustrated their performance with tapping data from the mPower study. We have also successfully applied our tests to features extracted from the voice and accelerometer data collected using the mPower app, but cannot present these results here due to space limitations.

About one quarter of the patients analyzed in this study showed strong statistical evidence for medication response. For these patients, we observed that the union-intersection tests tended to generate smaller p-values than the classifier based tests. This result corroborates our observations in the empirical power simulation studies, where union-intersection tests tended to out-perform classifier tests when the number of features is small (recall that we employed only 24 features in our analyses).

Although our tests can detect medication responses, they do not explicitly determine the direction of the response, that is, they cannot differentiate a response in the expected direction from a paradoxical one. We point out, however, that their main utility is to help out the physician pinpoint patients in need of a change on their drug treatment. For instance, our tests are able to detect patients for which the drug has an effect in the expected direction

but is not well calibrated (so that its effect is wore off by the time the patient takes the medication), or patients showing paradoxical responses. Note that both cases flag a situation which requires an action by the physician (calibrating the medication dosage in the first case, and stopping/changing the medication in the second one). Therefore, even though our tests cannot distinguish between these cases they are still able to detect patients which can potentially benefit from a dose calibration or a change in medication.

Even though we restricted our attention to 4 classifiers, other classification algorithms could also be easily used for generating additional personalized classifier tests. In our experience, however, tree based approaches such as the random forest and extra trees classifiers tend to perform remarkably well in practice, providing a robust default choice. Similarly, we focused on t- and Wilcoxon rank-sum tests in the derivation of our personalized union-intersection tests, since these tests show robust practical performance for detecting group differences.

Although others have investigated the feasibility of monitoring medication response in Parkinson patients using smartphone sensor data,⁴ their study focused on the classification of the before/after medication response at the population level, with the classifier applied to data across all patients, and not at a personalized level, as done here.

In this paper we show that simple hypothesis tests, which ignore the serial correlation in the feature data, are statistically valid and well powered to detect medication response in the mPower study data. We point out, however, that, in theory, more sophisticated approaches able to leverage the longitudinal nature of the data could in principle improve the statistical power to detect medication response. In any case, the good practical performance of our simple personalized tests suggests that they can be used as sound baseline approaches, against which more sophisticated methods can be compared.

Code and data availability. All the R¹³ code implementing the personalized tests, and used to generate the results and figures in this paper is available at https://github.com/Sage-Bionetworks/personalized_hypothesis_tests. The processed tapping data is available at doi:10.7303/syn4649804.

Acknowledgements. This work was funded by the Robert Wood Johnson Foundation.

References

1. S. Arora, et al, *Parkinsonism and related disorders* **21**, 650-653 (2015).
2. Editorial, *Nature Biotechnology* **33**, 567 (2015).
3. S. H. Friend, *Science Translational Medicine* **7**, 297ed10 (2015).
4. A. Zhan, et al, *High frequency remote monitoring of Parkinson's disease via smartphone: platform overview and medication response detection* (manuscript under review) (2015).
5. B. Efron, R. J. Tibshirani, *An introduction to the bootstrap*, Chapter 15 (Chapman and Hall/CRC, 1993).
6. J. Galambos, *Encyclopedia of Statistical Sciences* **7**, 573-577 (1996).
7. Y. Benjamini, Y. Hochberg, *Journal of the Royal Statistical Society, Series B* **57**, 289-300 (1995).
8. D. Bamber, *Journal of Mathematical Psychology* **12**, 387-415 (1975).
9. S. L. Mason, N. E. Graham, *Quarterly Journal of the Royal Meteorological Society* **128**, 2145-2166 (2002).
10. L. Breiman, *Machine Learning* **45**, 5-32 (2001).
11. P. Geurts, D. Ernst, L. Wehenkel, *Machine Learning* **63**, 3-42 (2006).
12. J. H. Friedman, T. Hastie, R. Tibshirani, *Journal of Statistical Software* **33** (1), (2010).
13. R Core Team, *R Foundation for Statistical Computing* (Vienna, Austria. ISBN 3-900051-07-0, URL <http://www.R-project.org/> 2015).