

1 Title: Explainable AI enables clinical trial patient selection to retrospectively improve treatment  
2 effects in schizophrenia

3

4 Authors: Monika S. Mellem\*, Matt Kollada, Jane Tiller, Thomas Lauritzen

5

6 Affiliations: BlackThorn Therapeutics, 780 Brannan St., San Francisco, CA 94103

7

8 \* Corresponding author

9

10

11

12

13

14

15

16

17

18

19

20

21

22

## 1 Abstract

2 **Background:** Heterogeneity among patients' responses to treatment is prevalent in psychiatric  
3 disorders. Personalized medicine approaches – which involve parsing patients into subgroups  
4 better indicated for a particular treatment – could therefore improve patient outcomes and  
5 serve as a powerful tool in patient selection within clinical trials. Machine learning approaches  
6 can identify patient subgroups but are often not “explainable” due to the use of complex  
7 algorithms that do not mirror clinicians' natural decision-making processes.

8 **Methods:** Here we combine two analytical approaches – Personalized Advantage Index and  
9 Bayesian Rule Lists – to identify paliperidone-indicated schizophrenia patients in a way that  
10 emphasizes model explainability. We apply these approaches retrospectively to randomized,  
11 placebo-controlled clinical trial data to identify a paliperidone-indicated subgroup of  
12 schizophrenia patients who demonstrate a larger treatment effect (outcome on treatment  
13 superior than on placebo) than that of the full randomized sample as assessed with Cohen's d.  
14 For this study, the outcome corresponded to a reduction in the Positive and Negative Syndrome  
15 Scale (PANSS) total score which measures positive (e.g., hallucinations, delusions), negative  
16 (e.g., blunted affect, emotional withdrawal), and general psychopathological (e.g., disturbance  
17 of volition, uncooperativeness) symptoms in schizophrenia.

18 **Results:** Using our combined explainable AI approach to identify a subgroup more responsive to  
19 paliperidone than placebo, the treatment effect increased significantly over that of the full  
20 sample ( $p < 0.0001$  for a one-sample t-test comparing the full sample Cohen's  $d = 0.82$  and a  
21 generated distribution of subgroup Cohen's  $d$ 's with mean  $d = 1.22$ , std  $d = 0.09$ ). In addition, our  
22 modeling approach produces simple logical statements (*if-then-else*), termed a “rule list”, to

1 ease interpretability for clinicians. A majority of the rule lists generated from cross-validation  
2 found two general psychopathology symptoms, disturbance of volition and uncooperativeness,  
3 to predict membership in the paliperidone-indicated subgroup.

4 **Conclusions:** These results help to technically validate our explainable AI approach to patient  
5 selection for a clinical trial by identifying a subgroup with an improved treatment effect. With  
6 these data, the explainable rule lists also suggest that paliperidone may provide an improved  
7 therapeutic benefit for the treatment of schizophrenia patients with either of the symptoms of  
8 high disturbance of volition or high uncooperativeness.

9 **Trial Registration:** clinicaltrials.gov identifier: NCT 00083668; registered May 28, 2004

10 **Keywords:** machine learning, personalized medicine, explainability, patient selection,  
11 schizophrenia

12

## 13 1. Background

14 The primary goal in a placebo-controlled clinical trial testing the efficacy of an  
15 experimental medication is to show a treatment effect – that patients randomized to receive  
16 the medication have improved outcomes compared to those receiving placebo. Within  
17 psychiatry, there is heterogeneity in patients’ responses, however, with some not responding  
18 well or at all [e.g., 1; 2] which can weaken the overall response of the treatment-receiving  
19 group compared to placebo. Additionally, the placebo response is robust in psychiatric  
20 disorders [3] making assessments of treatment efficacy more difficult. A method termed  
21 Personalized Advantage Index (PAI) has been recently developed to uncover subgroups of  
22 patients, termed “treatment-indicated,” may be more responsive to a particular treatment than

1 placebo and that predictive modeling could lead to personalized medicine approaches for  
2 subtyping treatment-indicated patients [4]. In particular, this could also help improve patient  
3 selection for clinical trials of that medication to enrich for patients most likely to show a  
4 treatment effect.

5 One of the limitations in using PAI to improve patient selection for clinical trials is  
6 insufficient explainability in how the model makes its decisions, as explainability is a critical  
7 attribute for a clinician to consider using an algorithm for patient selection. Prior work in  
8 depression has used the PAI approach to identify the most predictive variables [5], but  
9 interpretability of the models for clinicians could be further improved as they would require  
10 interpretation of regression coefficients and do not suggest clear cutoffs for predictor variables.  
11 Here, we additionally used an approach inspired by explainable artificial intelligence (XAI), the  
12 Bayesian Rule Lists algorithm (BRL) [6, 7], to both help identify the most predictive variables  
13 and explain those predictions of treatment-indicated patients from PAI with simple *if-then-else*  
14 statements that better mirror a clinician's decision-making process by using Boolean criteria  
15 with clear cutoffs for predictor variables. The combined analytical approach of PAI and BRL was  
16 previously tested in depression and found to retrospectively identify a subgroup with improved  
17 treatment effect for the novel antagonist BTRX-246040 [8]. But it has yet to be tested in other  
18 psychiatric populations.

19 While improving patient selection of a clinical trial is one potential use, it is important to  
20 note the constraints on this PAI and BRL approach for this goal and some additional  
21 opportunities. This approach requires both baseline and post-treatment (or imputed post-  
22 treatment since patients often do not discontinue at random) measurements of patients

1 receiving a particular treatment in order to learn the baseline characteristics of a treatment-  
2 indicated subgroup for this treatment prior to enrollment, so it may not be appropriate for  
3 clinical trial patient selection when no similar trial has been performed. One opportunity is  
4 thus to use it to learn the optimal subgroup from a negative clinical trial (as in [8]) and re-  
5 launch a more targeted clinical trial of the same treatment using the algorithm to identify a  
6 treatment-indicated subgroup. Enrolling this subgroup could help increase the treatment effect  
7 size of the more targeted clinical trial. A second opportunity is in using this approach to  
8 develop a decision-making support system for clinicians to prescribe medications already on  
9 the market to a targeted subgroup likely to have increased treatment efficacy. The data  
10 requirements of the PAI and BRL approach could be satisfied in all these scenarios.

11 In this study, we present this combined PAI and BRL approach to patient selection for  
12 clinical trials using XAI and validation in schizophrenia patients through a retrospective analysis  
13 of a clinical trial. By showing that baseline features alone can be used to identify a subgroup of  
14 patients who demonstrate a larger average treatment effect, this approach opens up the  
15 possibility of identifying a targeted subgroup of patient prior to randomization to an arm and  
16 improving the clinical trial outcome by using this patient subgroup in particular.

17

## 18 2. Methods

### 19 2.1 Study

20 We analyzed data from a 6-week randomized, double-blind, placebo-controlled study  
21 evaluating the efficacy of extended-release paliperidone in the treatment of patients with  
22 schizophrenia (clinicaltrials.gov identifier: NCT 00083668). This trial was a success, and

1 paliperidone currently has FDA approval in this population. This study makes use of de-  
2 identified data made available via the YODA Project (<https://yoda.yale.edu/>; research proposal  
3 number 2019-4080) and was exempt from ethical oversight.

4 Patients were assessed for eligibility and phenotyped using standard clinical  
5 assessments at baseline and randomized to one of several arms. In this analysis, we used data  
6 from the 15mg/day paliperidone arm (n=113) and the placebo arm (n=120). We selected the  
7 15mg/day arm over arms testing lower doses (3mg or 9mg per day) as it gave the greatest  
8 efficacy effect compared to placebo. The efficacy endpoint was the Positive and Negative  
9 Syndrome Scale (PANSS), a set of 30 questions administered by a trained clinician and scored on  
10 a 1-7 ordinal scale (7 is most severe). The PANSS was administered weekly. After dropping  
11 patients with missing baseline values used as features in the modeling (see below), 95 patients  
12 remained in the treatment arm and 102 in the placebo arm.

13

## 14 2.2 Modeling

15 Our approach to identify treatment-indicated patients and improve the explainability of  
16 the machine learning algorithm output classifying these patients involved combining two  
17 approaches. First, we used the Personalized Advantage Index (PAI) algorithm [4] to create an  
18 index/score that is then used to label a patient as treatment-indicated or rest-indicated (the  
19 rest of the subjects who are not treatment-indicated). This machine learning approach relies  
20 on multiple linear regression which provides some level of explainability through the  
21 coefficients of the predictors but is likely more complex than clinicians' decision making  
22 processes which are closer in form to decision trees or lists. Thus, in a second step, we used the

1 Bayesian Rule Lists classifier [6, 7] along with the predicted treatment- and rest-indicated labels  
2 from PAI to create a more explainable classifier using decision lists. The overview of the  
3 workflow is presented in Figure 1.

## 5 2.2.1 Personalized Advantage Index Modeling and Labeling

6 In our modeling approach, the PAI algorithm was first used to identify treatment-  
7 indicated patients. The PAI approach predicts actual and counterfactual outcomes (i.e., a  
8 patient's outcome for their assigned arm, drug or placebo, and the non-assigned arm) and  
9 calculates the difference between these two scores (as previously described in [4, 5]).

10 Briefly, we used an Elastic Net regressor (implemented in the python package scikit-  
11 learn), with a grid search for hyperparameter optimization across the range of  $\alpha = [0.001,$   
12  $0.01, 0.1, 1, 10]$  and  $l1\_ratio = [0.1, 0.5, 0.9]$ . The input features are listed in Table 1 where  
13 "baseline" refers to week 0 of the trial, which precedes treatment arm randomization. The  
14 outcome modeled was the 6-week post-treatment total PANSS score. Scores from patients  
15 missing this 6-week score were replaced with their last observation ( $n = 80$  patients) which was  
16 consistent with the approach used in the original clinical trial analysis. This multiple regression  
17 model then predicts the actual post-treatment PANSS score (on the patient's randomized arm)  
18 and the hypothetical counterfactual score (by substituting the other arm in the regression  
19 equation). PAI then returns a quantitative score for each patient that indicates the difference  
20 between these predicted drug and placebo outcomes with better performance on drug  
21 corresponding to a negative PAI score. A subsequent threshold then creates the two classes of  
22 treatment-indicated and rest-indicated with possible thresholds examined in descending steps

1 of 0.5 (0, -0.5, -1, -1.5,...). We selected the threshold that allocates ~50% of the sample to  
 2 treatment-indicated class to maintain a balanced data set for the BRL classifier. Please note  
 3 that selecting a PAI score threshold is a matter of balancing algorithmic needs (two-class  
 4 classifiers work best with balanced data), and clinical needs (clinicians may consider a specific  
 5 percent decrease in symptoms, such as at least 30%, a clinically-meaningful decrease), and  
 6 researchers who use this method in the future may need to reassess whether the 50% criteria  
 7 used here will work in their scenario. As PAI was responsible for generating the best possible  
 8 indication labels for BRL to train on (i.e., generating “ground truth” labels for BRL), the PAI  
 9 regression model was trained on all the data. Please note that these are relative ground-truth  
 10 labels – the best labels we can come up with but not perfect since the real ground-truth of  
 11 counterfactual predictions will never be known. Thus these labels *function* as ground-truth for  
 12 training BRL, as true ground-truth labels cannot be known with this clinical trial design. The grid  
 13 search of hyperparameters showed that  $\alpha = 0.1$ ,  $l1\_ratio = 0.1$  minimized the  $R^2$ .  
 14

Model Input Features	PAI Model Output	BRL Model Output
<u>Demographics</u> <ul style="list-style-type: none"> <li>• Age</li> <li>• Sex</li> </ul> <u>Symptom scales</u> <ul style="list-style-type: none"> <li>• 30 Baseline individual item PANSS scores</li> <li>• Baseline total PANSS score</li> <li>• Baseline daytime drowsiness and quality of sleep scores from the Sleep Visual Analog Scale (VAS)</li> <li>• Baseline Personal and</li> </ul>	<ul style="list-style-type: none"> <li>• Predicted actual and counterfactual week 6 PANSS scores</li> <li>• Numerical PAI score (predicted outcome on treatment – predicted outcome on placebo)</li> <li>• Labels: treatment-indicated, rest-indicated created from thresholded PAI score</li> </ul>	<ul style="list-style-type: none"> <li>• Labels: treatment-indicated and rest-indicated</li> </ul>



<p>Social Performance Scale score</p> <ul style="list-style-type: none"> <li>• Baseline total Schizophrenia Quality of Life Scale score</li> </ul> <p><u>Others</u></p> <ul style="list-style-type: none"> <li>• Randomized arm (treatment or placebo; PAI only)</li> <li>• Interactions of randomized arm with the other features (PAI only).</li> </ul>		
---	--	--

1 **Table 1:** Model inputs and outputs. Several sequential outputs are generated during PAI  
2 modeling and are listed in order of generation.

3

#### 4 2.2.2 BRL Modeling and Labeling for Additional Explainability

5 The BRL algorithm was used to create a more explainable model from the initial PAI  
6 treatment-indicated and rest-indicated results. BRL uses sequenced logical rules and Bayesian  
7 inference to make classifications [6, 7]. Here, we took in the same baseline features other than  
8 the randomized arm and interactions (Table 1) and classified patients as treatment- or rest-  
9 indicated using the BRL-generated *if-then-else* statements. A Bayesian rule list is composed of  
10 Boolean statements that evaluate if features fit certain criteria such as “If depression symptom  
11 score > 10” and the subsequent classification if the statement is true – “then, patient is  
12 treatment-indicated.” These statements are closer to a physician’s decision-making process  
13 than are the PAI regression outputs (feature coefficients without clear cutoffs for feature  
14 values). Hyperparameters were set at 3 for the max rule length (number of Boolean statements  
15 combined in an individual rule), 2 for the Bayesian prior hyperparameters of expected

1 individual rule length and 2 for the expected rule list length (excluding the final base case). We  
2 used a 5-fold cross-validation framework that generated a model (a rules list) on the 80% of  
3 training data and used it to classify the patients in the remaining 20% of test data as treatment-  
4 or rest-indicated. Thus five rules lists were generated from the five folds of cross-validation. As  
5 an additional output of this study, we generated a “final” BRL model by training the model on  
6 the full data set. While we have not tested this final BRL patient selection tool here on an  
7 external data set, others could use it if they have the appropriate data.

8

### 9 2.2.3 Comparing Treatment Effects

10 After classifying each patient as treatment- and rest-indicated labels using BRL, we  
11 assessed if treatment-indicated patients showed an improved treatment effect compared with  
12 the full sample. It is important to note that as the treatment-indicated patients were  
13 randomized to both drug and placebo arms, we were able to evaluate their actual post-  
14 treatment outcomes and to calculate their group-level treatment effect. Here, the treatment  
15 effect was assessed on the actual week 6 PANSS total scores of individual patients grouped by  
16 their treatment arm (using Cohen’s  $d$  as a measure of treatment effect size). After determining  
17 the labels from the BRL outputs, the actual week 6 PANSS total scores for the treatment-  
18 indicated patient subgroup were used to calculate the treatment effect for that subgroup.  
19 Then the treatment effect of the treatment-indicated subgroup was compared with the  
20 treatment effect of the full randomized sample. Our null hypothesis is that patient selection  
21 with BRL provides no improvement of the treatment effect, while our alternative hypothesis is  
22 that BRL does improve the treatment effect relative to that of the full sample. Thus to test this

1 statistically, we generated a distribution treatment effects by performing BRL 100 times and  
2 calculating the Cohen's d for each BRL-labeled treatment-indicated subgroup. We then  
3 compared this distribution of 100 Cohen's d's with the Cohen's d of the full sample using a two-  
4 sided, 1-sample t-test. Additionally, we assessed the consistency of classification for each  
5 subject across the five rule lists generated from the five folds of cross-validation.

6

### 7 3. Results

8 Our approach to identify treatment-indicated patients and improve the explainability of  
9 the machine learning algorithm output classifying these patients involved combining two  
10 approaches – PAI and BRL. While the final output and results are the treatment-labeled patient  
11 subgroup and the *if-then-else* rules list from the BRL model, the initial modeling with the PAI  
12 algorithm produced some interim results that we first examined. The PAI regression equation  
13 modeled actual week 6 PANSS scores using actual treatment arm assignment and several  
14 demographic and baseline symptom severity predictor variables (see Table 1). Figure 2 shows  
15 the comparison of measured week 6 PANSS scores with predicted scores (a perfect model  
16 would show all samples sitting on the  $x=y$  line). The  $R^2$  of the model shows it explained 58% of  
17 the variance (adjusted  $R^2 = 0.32$ ). Then, by substituting in the counterfactual randomization  
18 arm, the regression equation was used to make predictions of the week 6 PANSS scores if  
19 patients were receiving the counterfactual treatment. The difference between actual and  
20 counterfactual predictions were used to calculate the PAI scores (predicted score on treatment  
21 – predicted score on placebo) and determine the indication labels to be used by BRL. The  
22 distribution of PAI scores are shown in Figure 3, and a threshold of -9.5 generated two balanced

1 classes: treatment-indicated (n=100) and rest-indicated (n=97). For the treatment-indicated  
2 subgroup, this PAI threshold corresponded to a 30% reduction in average post-treatment  
3 PANSS scores for patients in the treatment arm relative to the scores of patients in the placebo  
4 arm and an improved treatment effect size (Cohen's  $d = 1.51$  relative to the full sample  $d =$   
5  $0.82$ ). Note that this treatment effect size is just an intermediate calculation as PAI is not the  
6 only step in our approach given that it does not provide the level of explainability that BRL  
7 does. An evaluation of only the PAI step's out of sample generalization and cross-validated  
8 treatment effect improvement is provided in the Supplementary Materials under the PAI  
9 Validation section.

10 After determining the indication labels from the thresholded PAI scores, we trained the  
11 BRL classifier on these labels using 5-fold cross-validation. For the training data, the BRL  
12 classifier had an average accuracy of 77.9% (standard deviation = 2.0%), an average Area Under  
13 the ROC Curve (AUC) of 0.83 (std = 0.02), and an average F1 score of 0.77 (std = 0.02) across the  
14 five folds. For the test data, the BRL classifier had an average accuracy of 74.1% (standard  
15 deviation = 5.1%) an average AUC of 0.76 (std = 0.04), and an average F1 score of 0.73 (std =  
16 0.04) across the five folds. Of the 197 patients from the full randomized sample, 87 were  
17 labeled treatment-indicated by the BRL algorithm, and the full confusion matrix of cross-  
18 validated labels are shown in Figure 4. On the full cross-validated test results, the overall  
19 accuracy was 74.1%, the AUC score = 0.74, and the F1 score = 0.73.

20 Comparison of the two arms (treatment, placebo) for the full sample v. the treatment-  
21 indicated group shows an increase in the Cohen's  $d$  between arms from 0.82 to 1.24 (Figure 5).  
22 We also assessed if this large increase in effect size was consistent for the BRL-classified

1 treatment-indicated subgroup and if it was significantly greater than the full sample effect size.  
2 We generated a distribution of Cohen's d's by performing the same BRL process 100 times and  
3 calculating the treatment effect for the treatment-indicated subgroup each time. We found  
4 that the treatment effects from these 100 iterations of subgrouping were statistically greater  
5 than the full sample treatment effect (BRL-classified treatment-indicated subgroup Cohen's d's  
6 mean = 1.22, std = 0.09; two-sided, 1-sample t-test:  $t=43.2$ ,  $p<0.0001$ ).

7 The five rule lists returned by the BRL classifier differ slightly across the five folds but did  
8 identify high disturbance of volition and high uncooperativeness as commonly identifying  
9 baseline features of patients who were more likely to respond on treatment than on placebo.  
10 Table 2 displays the five rule lists.

11 < PLACE TABLE 2 APPROXIMATELY HERE >

12 Though the rule lists show some differences across folds, we found that they still  
13 classified patients similarly when applied to the whole data set (not just the test set). To  
14 quantitatively assess the consistency of rule list classification, we classified each patient five  
15 times with the five rule lists. This gave us five labels (either treatment-indicated or rest-  
16 indicated) for each patient from which we can calculate the number of times that a patient is  
17 classified as treatment-indicated (max possible is five times). Figure 6 shows a histogram of the  
18 number of times that patients were labeled treatment-indicated. For consistent classifiers,  
19 most patients should be labeled treatment-indicated either five times or zero times (which  
20 corresponds to a patient labeled rest-indicated five times), and the labeling reflects this well  
21 according to our histogram. Additionally, most patients labeled by the BRL algorithm  
22 treatment-indicated four or five times were treatment-indicated according to the "ground

1 truth” PAI labels. As expected, this was reversed with most patients labeled by the BRL  
2 algorithm treatment-indicated zero times were rest-indicated according to the “ground truth”  
3 PAI labels. Thus, the five rule lists mostly subtype patients similarly though their wording can  
4 differ.

5 We additionally generated a final BRL model that could be validated on external data  
6 sets (Figure 7). This model generated a single rule list as it is trained on the full data set as  
7 opposed to the cross-validation approach which generated five rules lists across the five folds.  
8 High baseline uncooperativeness and high baseline disturbance of volition each remain  
9 predictive of treatment-indicated subgroup membership. For training and testing on the full  
10 data set, accuracy was 76%, AUC was 0.81, and the F1 score was 0.75.

11

#### 12 4. Discussion

13 With this retrospective analysis, we have technically-validated an approach of using a  
14 combination of machine learning methods to identify clinically-explainable rules that effectively  
15 subtype paliperidone-indicated schizophrenia patients who show an improved treatment effect  
16 over the full randomized sample. This extends the prior validation that demonstrated the  
17 method’s effectiveness in a clinical trial of a novel depression treatment [8].

18 While this validation was performed on a successful trial where the full sample already  
19 displayed the success of the experimental drug, we demonstrated than the treatment effect  
20 can be further improved with a patient selection approach. Statistically speaking, increasing  
21 the effect size can help decrease the enrollment numbers for patients, thereby possibly  
22 decreasing the cost and time of a new clinical trial (for example when running a more targeted

1 confirmatory target phase two trial after gathering data in a traditional phase two trial). Thus,  
2 in addition to patient selection that improves the treatment effect for unsuccessful trials [5, 8],  
3 our results suggest that patient selection could help clinical development even for treatments  
4 with stronger effects.

5 While the primary goal in this study was to validate a patient selection approach, the  
6 methodology also allows us to better understand a potential paliperidone-responsive subtype  
7 of schizophrenia. A majority of the rule lists generated from the cross-validation framework  
8 found two general psychopathology symptoms, disturbance of volition and uncooperativeness,  
9 to be predictive of membership in the paliperidone-indicated subgroup suggesting that  
10 paliperidone may be indicated for the treatment of patients with either of these symptoms.  
11 While subtyping schizophrenia is an active area of research [e.g., 2, 9], the higher severity of  
12 uncooperativeness and disturbance of volition seen in this paliperidone-indicated subgroup has  
13 not been previously described and should be externally validated. For this reason, we have  
14 included a single BRL rule list as a “final” model which could be tested by others interested in a  
15 paliperidone-indicated subgroup.

16 This approach is particularly useful in the context of selecting patients for clinical trial  
17 enrollment as the clinical trial outcome is dependent on large effects that are seen for patients  
18 in the treatment arm but not placebo arm. However, the proposed approach with the  
19 additional clinician-friendly explainability of BRL could make it more broadly useable as a  
20 clinical decision support system which are not commonly incorporating machine learning yet  
21 [10, 11]. The framework could be extended to accommodate multiple classes for indications of

1 multiple treatments. With the proper validation, this could provide clinicians with a tool to  
2 match the best of several possible treatments to a particular patient.

3           Some limitations with this approach remain. The PAI model does have some bias in its  
4 predictions as shown by the greater differences in measured and predicted week 6 PANSS  
5 scores in the larger and smaller ranges. This reflects a model that is underfit and could be due  
6 to missing predictor variables or due to using a linear rather than non-linear model. Future  
7 iterations of this approach could test using a non-linear approach such as generalized random  
8 forest [12] to improve predictions in this first step. As the shortcomings of the PAI step can  
9 affect model accuracy of the BRL step, there is additional incentive to improve predictions from  
10 the PAI step. Even with this weakness, the PAI model still provided adequate predictions to  
11 allow the BRL model to find a treatment-indicated subgroup with improved treatment effect  
12 size. Another limitation is the testing only within a single data set. A more robust approach  
13 would be to test the BRL model in a separate data set to assess generalization of the model and  
14 the proposed paliperidone-indicated schizophrenia phenotype. Some may question using the  
15 baseline total PANSS score as a predictor variable either due to its possible collinearity with  
16 other PANSS item scores or that the resulting use of lower PANSS scores to classify rest-  
17 indicated patients (therefore higher PANSS score is indirectly predictive of treatment-indicated  
18 patients) may be reflecting an effect of regression to the mean for the treatment-indicated  
19 patients. Unpublished analyses in our lab did not find major differences in performance or  
20 predictive features whether including or not including this variable. Additionally, the critical  
21 result is not that PANSS scores are reduced for the treatment-indicated patients, but that they  
22 are reduced much more on treatment than on placebo. Thus selecting patients as rest-



1 indicated based on the baseline total PANSS score and the implications that has for selecting  
2 treatment-indicated patients does not have any bearing on the improved difference seen  
3 between arms of the treatment-indicated patients. And finally, this modeling approach cannot  
4 currently handle longitudinal independent or dependent variables, but extending the  
5 framework of this method to more data types could expand its useability.

6

## 7 5. Conclusions

8 These results help to technically validate our explainable AI approach to patient selection for a  
9 clinical trial by identifying a subgroup of schizophrenia with an improved treatment effect.  
10 Importantly, this approach opens up the possibility of identifying a targeted subgroup of patient  
11 prior to randomization to an arm and improving the clinical trial outcome by using this patient  
12 subgroup in particular.

13

14

## 15 Abbreviations

16

17 AI: Artificial Intelligence, AUC: Area Under the ROC Curve, BRL: Bayesian Rule List, FDA: U.S.

18 Food and Drug Administration, PANSS: Positive and Negative Syndrome Scale , PAI:

19 Personalized Advantage Index, PBO: Placebo, ROC: Receiver Operating Characteristic, Rest-ind:

20 Rest-indicated, TRT: Treatment, TRT-ind: treatment-indicated, XAI: Explainable Artificial

21 Intelligence, YODA: Yale University Open Data Access

22

1    Declarations

2

3    Ethics approval and consent to participate

4

5    This study makes use of de-identified data made available via the YODA Project

6    (<https://yoda.yale.edu/>; research proposal number 2019-4080) and was exempt from ethical

7    oversight.

8

9    Consent for publication

10

11   Not applicable.

12

13   Availability of data and materials

14   The data that support the findings of this study are available from the YODA Project

15   (<https://yoda.yale.edu/>) but restrictions apply to the availability of these data, which were used

16   under license for the current study, and so are not publicly available.

17

18   Competing interests

19

20   All authors are current or previous employees of BlackThorn Therapeutics and hold stock or

21   stock options in BlackThorn Therapeutics.

1

## 2 Funding

3

4 Funding was provided by BlackThorn Therapeutics whose employees were involved designing  
5 this post-hoc analysis, performing the analysis, interpreting the results, and writing the  
6 manuscript.

7

## 8 Authors' contributions

9

10 JT, MK, MM, and TL helped conceive of the approach to PAI+BRL validation with the YODA data  
11 set and subsequent analysis, and JT advised on methodology with the team who developed the  
12 PAI+BRL approach. MM performed the analysis of the data. MM, MK, and TL interpreted the  
13 results, and wrote or revised the manuscript. All authors read and approved the final  
14 manuscript.

15

## 16 Acknowledgments

17

18 We would like to thank Clark Gao, Yuelu Liu, Humberto Gonzalez, and Parvez Ahammad for  
19 their contributions in creating the PAI+BRL methodology that is used here in a different  
20 application, and Kathy Tracy, Rezi Zawadzki, Andrew Jaffe, John Dunlop, and Bill Martin for  
21 further discussions. This study, carried out under YODA Project # 2019-4080, used data  
22 obtained from the Yale University Open Data Access Project, which has an agreement with

1 JANSSEN RESEARCH & DEVELOPMENT, L.L.C. The interpretation and reporting of research using  
2 this data are solely the responsibility of the authors and does not necessarily represent the  
3 official views of the Yale University Open Data Access Project or JANSSEN RESEARCH &  
4 DEVELOPMENT, L.L.C.

5  
6

## 7 References

8

- 9 1. Akil H, Gordon J, Hen R, Javitch J, Mayberg H, McEwen B, Meaney MJ, Nestler EJ.  
10 Treatment resistant depression: a multi-scale, systems biology approach. *Neuroscience*  
11 *and Biobehavioral Reviews*. 2018; 84: 272-288.  
12 [doi.org/10.1016/j.neubiorev.2017.08.019](https://doi.org/10.1016/j.neubiorev.2017.08.019)  
13
- 14 2. Gillespie AL, Samanaite R, Mill J, Egerton A, MacCabe JH. Is treatment-resistant  
15 schizophrenia categorically distinct from treatment-responsive schizophrenia? a  
16 systematic review. *BMC Psychiatry*. 2017; 17(12):1-14. [doi.org/10.1186/s12888-016-](https://doi.org/10.1186/s12888-016-1177-y)  
17 [1177-y](https://doi.org/10.1186/s12888-016-1177-y)  
18
- 19 3. Weimer K, Colloca L, Enck P. Placebo effects in psychiatry: mediators and moderators.  
20 *Lancet Psychiatry*. 2015; 2(3):246-257. [doi:10.1016/S2215-0366\(14\)00092-3](https://doi.org/10.1016/S2215-0366(14)00092-3)  
21

- 1        4. DeRubeis RJ, Cohen ZD, Forand NR, Fournier JC, Gelfand LA, Lorenzo-Luaces L. The  
2            personalized advantage index: translating research on prediction into individualized  
3            treatment recommendations. A demonstration. PLoS One. 2014; 9(1):e83875.  
4            doi:10.1371/journal.pone.0083875  
5
- 6        5. Webb CA, Trivedi MH, Cohen ZD, Dillon DG, Fournier JC, Goer F, et al. Personalized  
7            prediction of antidepressant v. placebo response: evidence from the EMBARC study.  
8            Psychological Medicine. 2018; 1-10. doi.org/10.1017/S0033291718001708  
9
- 10       6. Gao Q, Gonzalez H, Ahammad P. MCA-based rule mining enables interpretable inference  
11           in clinical psychiatry. In Shaban-Nejad A., Michalowski M. (eds) Precision Health and  
12           Medicine. W3PHAI 2019. Studies in Computational Intelligence. Springer, Cham. 2020;  
13           vol. 843. doi.org/10.1007/978-3-030-24409-5\_3  
14
- 15       7. Letham B, Rudin C, McCormick TH, Madigan D. Interpretable classifiers using rules and  
16           Bayesian analysis: building a better stroke prediction model. The Annals of Applied  
17           Statistics 2015; 9(3):1350-1371. doi.org/10.1214/15-AOAS848.  
18
- 19       8. Martin B, Gao Q, Liu Y, Madrid A, Ahammad P, Tiller J. Explainable AI approach reveals  
20           treatment responders in a randomized controlled trial of BTRX-246040, a potent and  
21           selective NOP receptor antagonist. Neuropsychopharmacology, 2019, vol. 44, no. SUPPL  
22           1, pp. 448-448.  
23

1 9. Jablensky A. Subtyping schizophrenia: implications for genetic research. Mol. Psychiatry.  
2 2006; 11(9):815-836. doi: 10.1038/sj.mp.4001857.

3  
4 10. Levy-Fix G, Kuperman GJ, Elhadad N. Machine learning and visualization in clinical  
5 decision support: current state and future directions. 2019; Preprint at  
6 [arXiv:1906.02664v1](https://arxiv.org/abs/1906.02664v1)

7  
8 11. Sutton RT, Pincock D, Baumgart DC, Sadowski DC, Fedorak RN, Kroeker KI. An overview  
9 of clinical decision support systems: benefits, risks, and strategies for success. NPJ  
10 Digital Medicine. 2020; 3(17). doi.org/10.1038/s41746-020-0221-y

11  
12 12. Athey S, Tibshirani J, Wager S. Generalized random forests. Ann. Statist. 2019;  
13 47 (2) 1148 - 1178. <https://doi.org/10.1214/18-AOS1709>

14

## 15 Figure Legends

16

17 **Figure 1:** Overview of PAI and BRL modeling workflow. **A)** The first step is PAI regression  
18 modeling which takes in data listed in Table 1 and trains on the whole data set to predict both  
19 actual and counterfactual post-treatment scores for individual patients (actual scores can be  
20 compared with predicted actual scores and resulted in an  $R^2 = 0.58$  as shown in Figure 2). **B)** The  
21 PAI Thresholding step thresholds the difference between actual and counterfactuals to create  
22 indication labels for each patient. A treatment-indicated subgroup had a treatment effect size

1 of Cohen's  $d = 1.51$  as an intermediate assessment, but explainability of model decisions needs  
2 improvement, so the BRL step addresses this need. **C)** The BRL modeling uses 5-fold cross-  
3 validation to assess generalization ability to unseen samples. The predictions generated for test  
4 samples over all folds had an accuracy of 74.1% and an AUC of 0.74 for this classifier.  
5 Importantly, it emits an explainable rule list for each fold. **D)** The final step is assessing the  
6 treatment effect of the treatment-indicated subgroup identified by BRL (Cohen's  $d = 1.24$  as  
7 seen in Figure 5).

8  
9 **Figure 2:** Individual PAI prediction results for actual randomized arms. The plot shows the  
10 measured total PANSS score at week 6 vs. the averaged predicted total PANSS score at week 6  
11 from the Elastic Net regression model (each dot is an individual patient,  $n=197$ ). Patients are  
12 colored by their actual randomized arm (paliperidone treatment in blue, placebo in orange), and  
13 as expected the week 6 scores are generally higher for patients receiving placebo. The dashed  
14 line is  $y=x$ . Variance explained by the model is 58%. Note that these are not the counterfactual  
15 predictions.

16  
17 **Figure 3:** PAI score thresholded graph. A threshold of -9.5 was chosen to create roughly  
18 balanced classes and indicates that membership in the treatment-indicated subgroup required a  
19 predicted treatment arm PANSS score that is 9.5 points less than the predicted placebo arm  
20 PANSS score.

21

1 **Figure 4:** Confusion matrix for cross-validated BRL labels. Actual labels are PAI-derived labels,  
2 and predicted labels are BRL-derived labels.

3  
4 **Figure 5:** Comparison of actual post-treatment PANSS scores for the full sample and the  
5 treatment-indicated subgroup. Bars display treatment (TRT) and placebo (PBO) arms for the full  
6 randomized sample (left graph, TRT n=95, PBO n=102) with an illustrative instance of the BRL-  
7 classified treatment-indicated (TRT-indicated) subgroup (right graph, TRT n=41, PBO n=46). At a  
8 Cohen's d of 1.24, the effect size between arms for the treatment-indicated subgroup is  
9 increased more than 50% over the effect size of the full sample ( $d = 0.82$ ). Error bars are 95%  
10 confidence intervals on the mean.

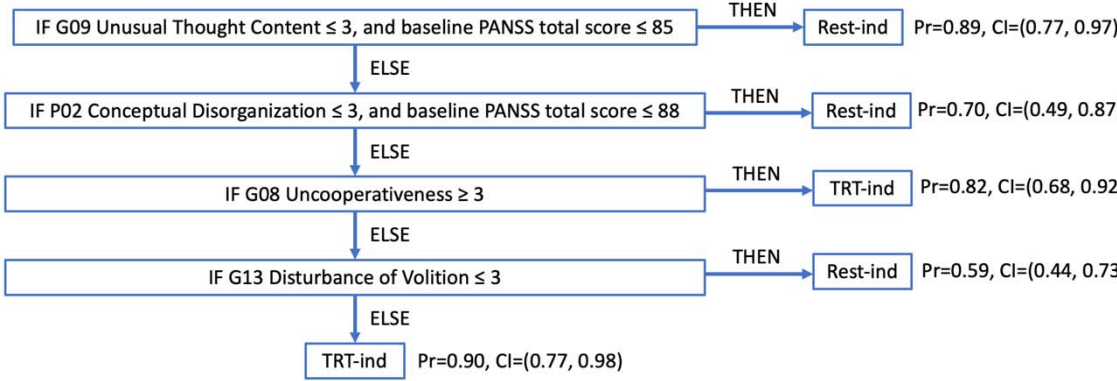
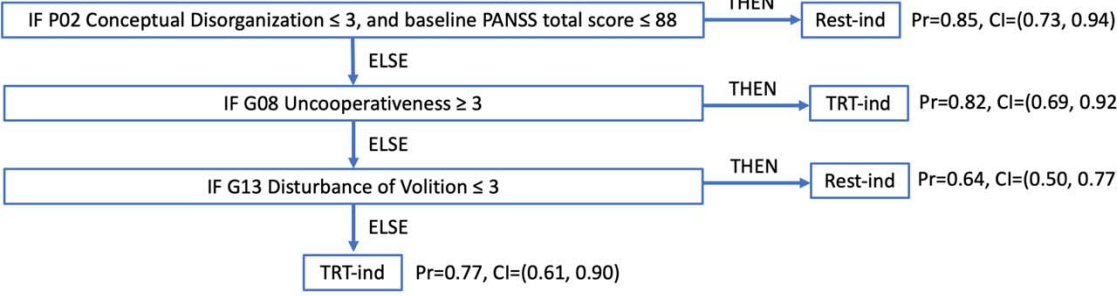
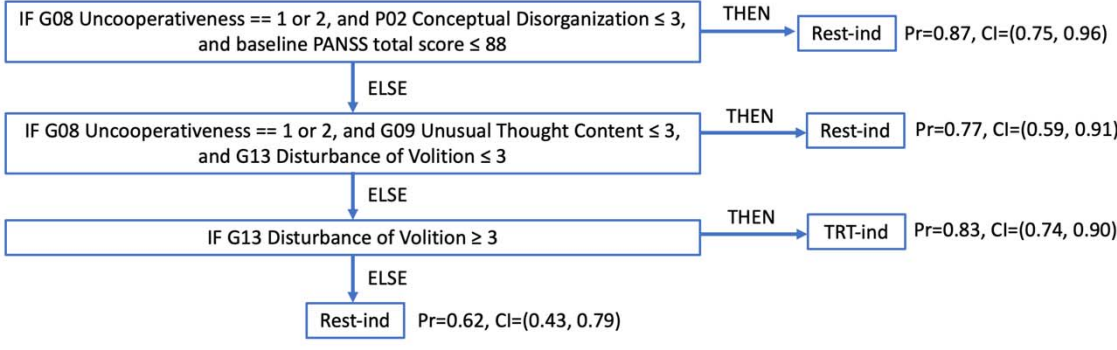
11  
12 **Figure 6:** The treatment-indicated (TRT-ind) labeling consistency is seen with a histogram of  
13 BRL-labeled treatment-indicated counts. It reflects the number of times (count) that a patient  
14 was classified as treatment-indicated across all five rule lists from the 5-fold BRL cross-  
15 validation. Most patients are either classified five times or zero times which indicates a higher  
16 level of consistency in patient subtyping across the different rule lists. Additionally, most patients  
17 in the 5-count column were also labeled as TRT-ind by the PAI algorithm (orange portion of the  
18 bar) while most patients in the zero-count column were labeled as Rest-ind by PAI (blue portion  
19 of the bar). Patient numbers corresponding to the orange and blue portions are shown in the  
20 table below the histogram.

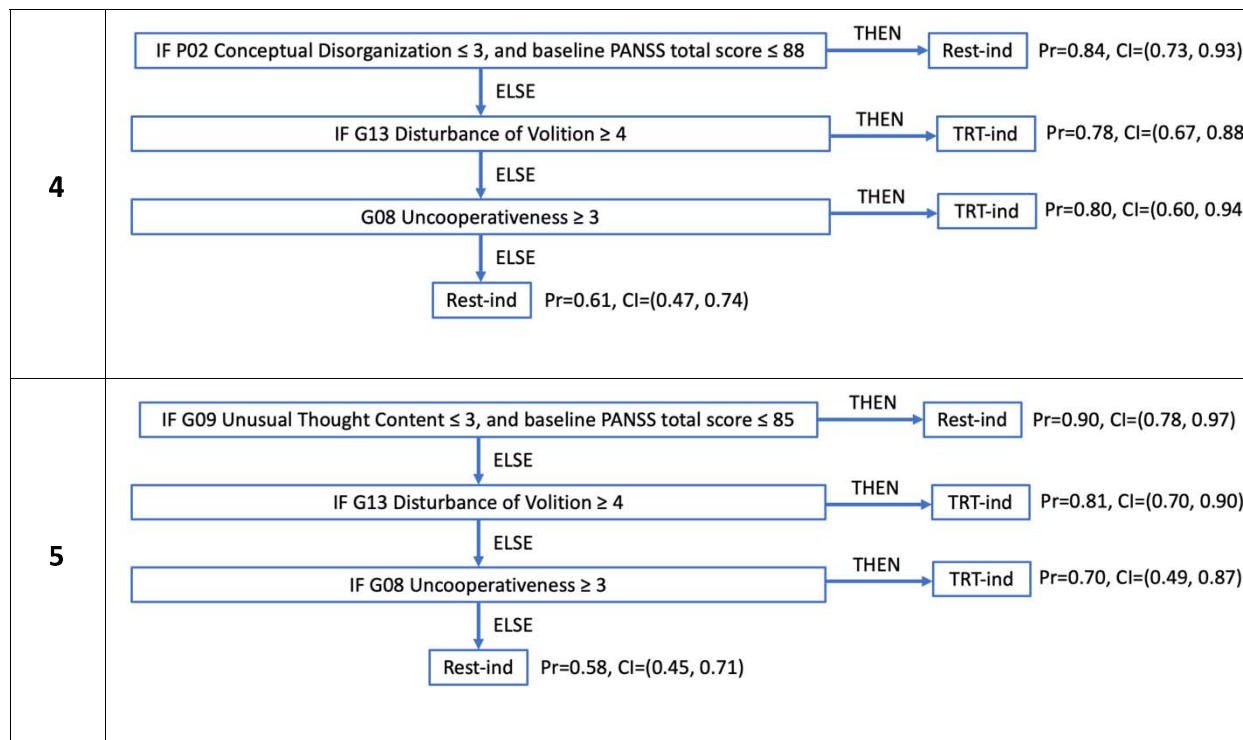
21



1 **Figure 7:** “Final” rule list created by the BRL model when trained on the full data set. The two  
 2 individual item scores for Disturbance of Volition and Uncooperativeness again appear to  
 3 classify the treatment-indicated (TRT-ind) subgroup.

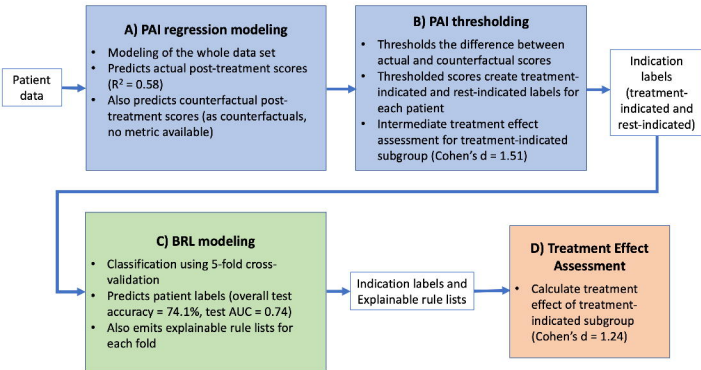
4  
 5 **Table 2:** List of the five rule lists created by the BRL model across the five folds.

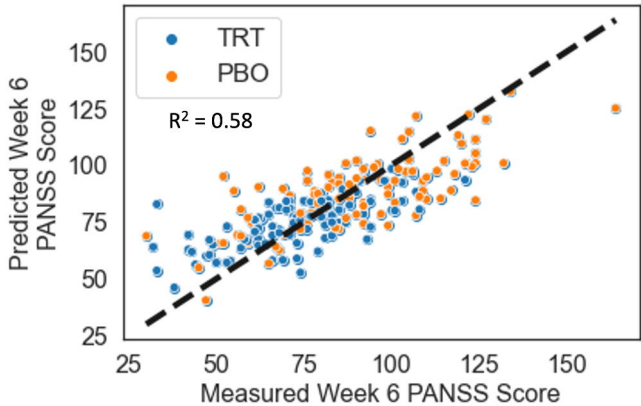
Fold	Rule List
1	 <pre> graph TD     A["IF G09 Unusual Thought Content ≤ 3, and baseline PANSS total score ≤ 85"] -- THEN --&gt; B["Rest-ind Pr=0.89, CI=(0.77, 0.97)"]     A -- ELSE --&gt; C["IF P02 Conceptual Disorganization ≤ 3, and baseline PANSS total score ≤ 88"]     C -- THEN --&gt; D["Rest-ind Pr=0.70, CI=(0.49, 0.87)"]     C -- ELSE --&gt; E["IF G08 Uncooperativeness ≥ 3"]     E -- THEN --&gt; F["TRT-ind Pr=0.82, CI=(0.68, 0.92)"]     E -- ELSE --&gt; G["IF G13 Disturbance of Volition ≤ 3"]     G -- THEN --&gt; H["Rest-ind Pr=0.59, CI=(0.44, 0.73)"]     G -- ELSE --&gt; I["TRT-ind Pr=0.90, CI=(0.77, 0.98)"]     </pre>
2	 <pre> graph TD     A["IF P02 Conceptual Disorganization ≤ 3, and baseline PANSS total score ≤ 88"] -- THEN --&gt; B["Rest-ind Pr=0.85, CI=(0.73, 0.94)"]     A -- ELSE --&gt; C["IF G08 Uncooperativeness ≥ 3"]     C -- THEN --&gt; D["TRT-ind Pr=0.82, CI=(0.69, 0.92)"]     C -- ELSE --&gt; E["IF G13 Disturbance of Volition ≤ 3"]     E -- THEN --&gt; F["Rest-ind Pr=0.64, CI=(0.50, 0.77)"]     E -- ELSE --&gt; G["TRT-ind Pr=0.77, CI=(0.61, 0.90)"]     </pre>
3	 <pre> graph TD     A["IF G08 Uncooperativeness == 1 or 2, and P02 Conceptual Disorganization ≤ 3, and baseline PANSS total score ≤ 88"] -- THEN --&gt; B["Rest-ind Pr=0.87, CI=(0.75, 0.96)"]     A -- ELSE --&gt; C["IF G08 Uncooperativeness == 1 or 2, and G09 Unusual Thought Content ≤ 3, and G13 Disturbance of Volition ≤ 3"]     C -- THEN --&gt; D["Rest-ind Pr=0.77, CI=(0.59, 0.91)"]     C -- ELSE --&gt; E["IF G13 Disturbance of Volition ≥ 3"]     E -- THEN --&gt; F["TRT-ind Pr=0.83, CI=(0.74, 0.90)"]     E -- ELSE --&gt; G["Rest-ind Pr=0.62, CI=(0.43, 0.79)"]     </pre>

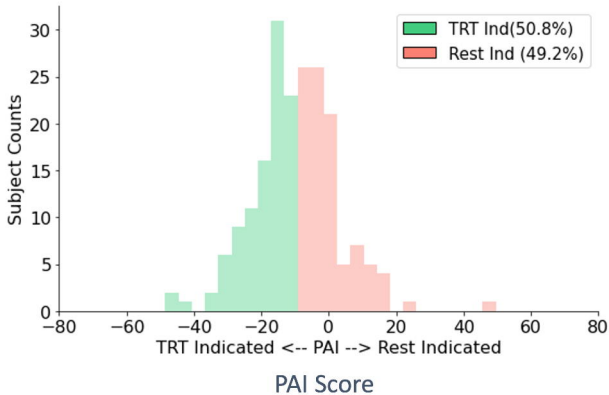


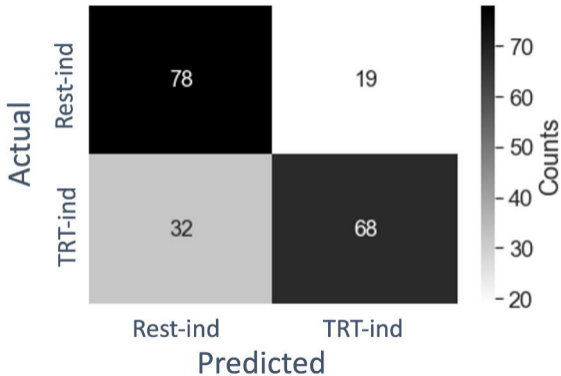
1 **Table 2 Legend:** The probability shown in parentheses after each rule is estimated from the  
 2 percent of patients who satisfy that rule and were labeled with the given label for the rule by  
 3 the PAI “ground-truth” labels, and the confidence intervals were estimated with bootstrapping.  
 4 The alphanumeric symbol before each symptom (e.g., P02, G08) refers to the question number  
 5 from the PANSS scale. Two individual item scores from the PANSS scale repeatedly were  
 6 involved in subtyping the treatment-indicated patients. Both Disturbance of Volition  $\geq 3$  and  
 7 Uncooperativeness  $\geq 3$  appear in multiple rule lists for the treatment-indicated (TRT-ind)  
 8 subgroup and are categorized as General Psychopathology symptoms.

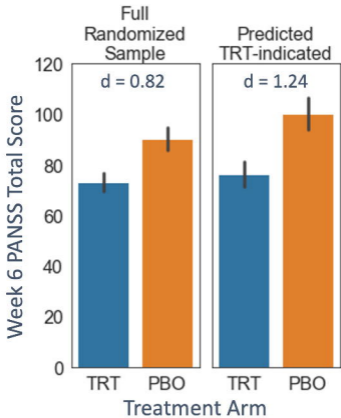
9

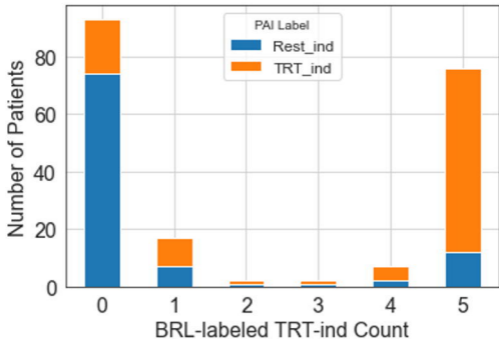












BRL Count	0	1	2	3	4	5
PAI TRT-ind	19	10	1	1	5	64
PAI Rest-ind	74	7	1	1	2	12



IF G09 Unusual Thought Content  $\leq 3$ , and baseline PANSS total score  $\leq 85$

THEN

Rest-ind

Pr=0.89, CI=(0.78, 0.96)

ELSE

IF G08 Uncooperativeness  $\geq 3$

THEN

TRT-ind

Pr=0.78, CI=(0.66, 0.88)

ELSE

IF G13 Disturbance of Volition  $\geq 4$

THEN

TRT-ind

Pr=0.81, CI=(0.67, 0.92)

ELSE

Rest-ind

Pr=0.61, CI=(0.49, 0.72)