medRxiv preprint doi: https://doi.org/10.1101/19010041; this version posted October 29, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a CC-BY 4.0 International license.

# **BRIEF REPORT**

31

32

33

34

35

36

37

38

39

40

41

# Machine learning-based imaging biomarkers improve statistical power in clinical trials

Carolyn Lou<sup>a,b</sup>, Mohamad Habes<sup>b</sup>, Christos Davatzikos<sup>b,1,\*</sup>, and Russell T. Shinohara<sup>a,b,1,\*</sup>

<sup>a</sup> Penn Statistics in Imaging and Visualization Center, Department of Biostatistics, Epidemiology, and Informatics, University of Pennsylvania, 423 Guardian Drive, Philadelphia, PA, 19104; <sup>b</sup>Center for Biomedical Image Computing and Analytics, Department of Radiology, University of Pennsylvania, 3700 Hamilton Walk, Philadelphia, PA, 19104; <sup>c</sup>Data used in preparation of this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database (adni.loni.usc.edu). As such, the investigators within the ADNI contributed to the design and implementation of ADNI and/or provided data but did not participate in analysis or writing of this report. A complete listing of ADNI investigators can be found at: http://adni.loni.usc.edu/wp-content/uploads/how\_to\_apply/ADNI\_Acknowledgement\_List.pdf.; <sup>\*</sup>These authors contributed equally.

This manuscript was compiled on October 18, 2019

Radiomic models, which leverage complex imaging patterns and ma-1 chine learning, are increasingly accurate in predicting patient re-2 sponse to treatment and clinical outcome on an individual patient 3 basis. In this work, we show that this predictive power can be utilized in clinical trials to significantly increase statistical power to de-5 tect treatment effects or reduce the sample size required to achieve 6 a given power. Akin to the historical control paradigm, we propose 7 to utilize a radiomic prediction model to generate a pseudo-control 8 sample for each individual in the trial of interest. We then incorpo-9 rate these pseudo-controls into the analysis of the clinical trial of 10 interest using classical and well established statistical tools, and in-11 vestigate statistical power. Effectively, this approach utilizes each 12 individual's radiomics-based predictor of outcome for comparison 13 with the actual outcome, potentially increasing statistical power con-14 siderably, depending on the accuracy of the predictor. In simulations 15 of treatment effects based on real radiomic predictive models from 16 brain cancer and prodromal Alzheimer's Disease, we show that this 17 methodology can decrease the required sample sizes by as much 18 as a half, depending on the strength of the radiomic predictor. We 19 further find that this method is most helpful when treatment effect 20 sizes are small and that power grows with the accuracy of radiomic 21 prediction. 22

Machine Learning | Personalized Medicine | Clinical Trials

n recent decades, rapid advances in technology have increased the amount of neuroimaging data available to re-2 searchers at an unprecedented rate (1, 2). Machine learning 3 methods empower the integration of this high-dimensional 4 data into powerful individualized predictive markers that have 5 been shown to be useful for tasks such as diagnosis and prog-6 nosis in diseases such as Alzheimer's disease and brain cancers (3, 4). Predictive modeling is poised to receive the benefits of 8 the large and varied nature of this data. 9

With the growing availability of big data in medical imaging, 10 a central focus has emerged on the development of increas-11 ingly complex tools for their analysis with the primary goal of 12 individualized predictions (5). In this paper, we propose har-13 nessing these powerful machine learning tools for the analysis 14 of clinical trials by using them as a means to inform statistical 15 analyses with individualized estimates of clinical outcome. We 16 therefore arrive at the concept of individualized evaluation of 17 treatment effects in clinical trials. 18

There is an extensive literature on the use of historical controls to supplement data from new clinical trials (6, 7). These methods have largely relied on pooling methods or Bayesian modeling. Whereas these methods augment data for a current trial by incorporating historical data on the 23 group level, high-dimensional predictors offer the opportunity 24 to augment current trials by incorporating historical data to 25 develop individualized predictions, or synthetic control data 26 (8), at the individual level. This allows for a more precise 27 evaluation of the treatment effect for each person, rather 28 than relying on a group-level effect that determines average 29 outcome. 30

Here, we present a method that draws on these ideas while leveraging powerful predictive biomarkers and the wealth of data used to build them to generate personalized predictions of outcome. These predictions can be used directly in the analysis of data in clinical trials. We find that this methodology can substantially improve statistical power for detecting treatment effects, depending on the predictive power of the machine learning-based model. Correspondingly, this approach can substantially reduce the sample size needed to achieve the same power in a clinical trial.

# Methods

Our method relies on access to two sets of data: i) a current 42 clinical trial designed to study an outcome of interest and ii) 43 a cohort of similar subjects treated according to the current 44 standard of care. We narrow our focus in this work to radiomic 45 predictors and associated studies, so we assume that imaging 46 data has been gathered at study enrollment for both sets of 47 trials. However, more broadly we only require a predictive 48 model that is based on sufficient information measured at 49 baseline on each participant in both datasets to predict the 50 outcome under standard of care. The techniques proposed here 51 are also directly applicable to other -omic modeling scenarios, 52 and generally, to any predictive marker of standard of care 53 outcome. 54

Our basic premise herein is that we can utilize previously 55 collected imaging data to build a radiomic prediction model, 56 fully validate it, and use it to generate a single score that 57 summarizes imaging patterns that predict future clinical out-58 come of interest, such as patient survival, progression-free 59 survival, or response to treatment (Figure 1). The model 60 that is built based on the historical trial can then be used 61 in conjunction with data collected from the current trial to 62

The authors declare no conflict of interest.

<sup>1</sup>To whom correspondence may be addressed. E-mails: christos.davatzikos@pennmedicine.upenn.edu or rshi@pennmedicine.upenn.edu

R.T.S. and C.D. designed research; C.L. performed research; C.L. and M.H. analyzed data; and C.L., M.H., R.T.S., and C.D. wrote the paper.

medRxiv preprint doi: https://doi.org/10.1101/19010041; this version posted October 29, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a CC-BY 4.0 International license .



Fig. 1. A: Workflow for implementing the proposed method in a new clinical trial. B and C: Schematic diagram for individualized predictions that are generated for each person in the current trial, where the solid lines indicate observed outcome for the participants of the current trial and the dashed lines indicate predicted outcome for those participants had they not been treated. Figure B illustrates the method for continuous outcomes, where the left side represents the controls and the right side represents the treated participants. The predicted outcome values (vertical axis) for the control units had they not been treated would be exactly what they are observed to be, while the predicted outcome values for the treated units had they not been treated are different from the observed outcome. Figure C illustrates the mechanism for survival outcomes, where the predicted survival times (dashed lines) for the control units are the same as the observed survival times (solid lines), whereas the predicted survival for the treated individuals are lower than the observed survival times.

generate individualized values of the radiomic score for each of 63 the current participants. These individualized scores represent 64 predicted values for how the treated individuals in the current 65 trial would have fared had they instead been assigned to the 66 control group. The incorporation of these predicted values 67 lends power to the detection of the effect of a treatment in the 68 final analysis of the current trial by modeling the inter-subject 69 variability in the outcome in terms of baseline heterogeneity 70 represented in the baseline imaging. 71

To investigate the advantage of this approach, we consider 72 two scenarios. To more closely approximate real-life clinical 73 trial performance, we use radiomic and outcome data from 74 two observational studies to generate hypothetical study data, 75 where the first focuses on continuous outcomes of cognitive 76 decline in prodromal Alzheimer's disease (AD) and the sec-77 ond on survival after diagnosis with glioblastoma multiforme 78 (GBM). In these studies, we randomly split the data into a 79 historical cohort and a trial cohort and then simulate effects 80 in a randomly selected subset (corresponding to one arm) of 81 the trial cohort. We then compare the statistical power of our 82 proposed approach with the classical modeling approach that 83 does not include radiomic prediction-based modeling. 84

Data. In our first study, we consider the case of therapeutic 85 trials for AD in which the outcome is longitudinal cognitive 86 change. We utilize data derived from the Azheimer's Disease 87 Neuroimaging Initiative (ADNI, adni.loni.usc.edu) on 400 sub-88 jects with mild cognitive impairment (MCI) who underwent 89 serial MRI at  $1.5T^*$ . 90

Data used from this study consisted of cross-validated pre-91 dictions of time to AD diagnosis using the SPARE-AD score 92 (9). SPARE-AD is derived from patterns of regional brain 93 atrophy (volume loss) captured by atlas warping methods and 94 high-dimensional pattern classification using support vector 95 96 machines (SVM) aiming to differentiate cognitively normal

and Alzheimer's disease subjects (3, 10). The outcome of interest here is cognitive decline as measured by 3-year change from baseline values of the ADNI composite memory score (11) (ADNI-MEM). Of the 400 MCI subjects in our study, 100 283 have 3-year ADNI-MEM scores available. Average 3-year 101 change from baseline for ADNI-MEM in the current study was 102 -0.17 (standard deviation 0.49).

As a second case study, we focus on the apeutic trials for 104 GBM therapies, which aim to prolong survival after diagnosis. 105 We conduct our simulated treatment study using 134 patients 106 who were treated for newly diagnosed GBM at the Hospital of 107 the University of Pennsylvania between 2006 and 2013. The 108 actual median survival in this sample was 12 months, and 109 survival data were assessed for all subjects with no loss to 110 follow-up. Detailed demographics and a clinical description of 111 these subjects have been previously published (4). For studies 112 involving these data, we investigate the use of cross-validated 113 predictions of survival time based on radiomic analyses of 114 pre- and post-contrast T1-weighted, FLAIR, diffusion, and 115 perfusion imaging acquired pre-operatively at diagnosis. This 116 GBM predictive model utilizes an SVM to differentiate short, 117 medium, and long survival (4). 118

Statistical Methods. All hypothesis testing is conducted assum-119 ing a 5% type I error rate and using two-sided alternatives. For 120 our continuous outcome analyses, we apply linear regression 121 modeling of the outcome and employ Wald tests to assess 122 whether treatment groups differed in their outcomes either i) 123 adjusting for the radiomic predictor by inclusion as covariate, 124 or the classical approach with corresponds to ii) not adjusting 125 for the radiomic predictor. For time-to-event outcomes, we 126 assess differences between treatment groups with and with-127 out adjustment for the radiomic prediction by assuming an 128 accelerated failure time model. 129

We conduct two sets of real data simulations: one set 130 focusing on cognitive decline in AD, and one set focusing 131 on GBM survival outcomes. For both, we sample without 132 replacement twice from the observed data: for the first group 133 indexed by  $i = 1, \ldots, n/2$ , we set our treatment indicator 134  $A_i = 0$  and record the observed outcome  $Y_{i0}$ , as well as the 135

97

98

99

103

<sup>\*</sup>The ADNI was launched in 2003 as a public-private partnership, led by Principal Investigator Michael W. Weiner, MD. The primary goal of ADNI has been to test whether serial magnetic resonance imaging (MRI), positron emission tomography (PET), other biological markers, and clinical and neuropsychological assessment can be combined to measure the progression of mild cognitive impairment (MCI) and early Alzheimer's disease (AD). For up-to-date information, see www.adni info.org

medRxiv preprint doi: https://doi.org/10.1101/19010041; this version posted October 29, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a CC-BY 4.0 International license.



Fig. 2. Results from simulated studies under two scenarios. With the addition of historical controls, the required sample size for 80% power is markedly lower than using classical two-sample clinical trial analysis. These figures show minimum sample size (vertical axes) required to achieve 80% power for a range of effect sizes (horizontal axes) based on observed outcome and radiomic predictions. Figure A shows the results from simulations for continuous outcome measures of cognition in AD, and Figure B shows the results from simulations for survival in GBM (right). Note that the proposed method that leverages historical controls to build radiomic predictions (red) requires lower samples sizes than the classical approach (blue).

value of the radiomic predictor  $X_i$  at baseline. For the second 136 group, indexed by  $i = n/2 + 1, \ldots, n$ , we introduce a treatment 137 effect  $\gamma$ , set our treatment indicator  $A_i = 1$ , and again record 138 outcome  $Y_{i0}$  and baseline radiomic predictor measurement 139 140  $X_i$ . We repeat this process 1000 times, recording the p-value 141 corresponding to the test for treatment effect each time. We calculate type I error rate and power as the percentage of time 142 the treatment effect is significant at the  $\alpha = 0.05$  level, where 143  $\gamma$  is set to 0 to assess type I error and a non-zero value to assess 144 power. In order to quantify the sample size benefits from using 145 this method, we repeat the above procedure for a range of 146 sample sizes n, and the smallest n for which power reaches 147 80% is recorded. We explore this for a range of hypothetical 148 effect sizes, which is defined here as  $\gamma$  divided by the standard 149 deviation of the outcome. 150

## 151 Results

For both continuous and time-to-event outcomes, we find that 152 the proposed method reduces the minimum sample size n153 required for 80% power in clinical trial analyses (Figure 2). 154 The inclusion of the imaging biomarker tends to be most 155 helpful in terms of the absolute differences when the effect 156 sizes are small. Type I error remains controlled throughout 157 all experiments conducted. In the ADNI study, the classical 158 159 analysis requires 16% to 18% more samples than our proposed method. In the GBM study, the classical analysis requires 160 around 73% to 94% more samples than our proposed method. 161

## 162 Discussion

We have shown that individualized machine learning-based 163 imaging biomarkers can be a useful tool in a clinical trial 164 analysis, offering an increase in power and/or a reduction of 165 the required sample size. The novelty of this method arises 166 from the incorporation of individualized predictions, which 167 derive their usefulness from the powerful predictive algorithms 168 they are based upon. As neuroimaging biomarkers derived via 169 machine learning become more common, the set of historical 170

data for which we have biomarker values also becomes larger, 171 which promises to strengthen the radiomic prediction model 172 that we use to generate predicted values. 173

Here, we used two previously developed biomarkers, one 174 of which was trained to classify an outcome different from 175 the target of the clinical trial analysis, and another which 176 was trained to classify the same outcome as the clinical trial 177 analysis. While both predictors offered gains in sample size 178 reduction, predictors built specifically for the outcome of in-179 terest in the clinical trial are likely to perform better and offer 180 more substantial gains. 181

The approach proposed in this paper does have some limi-182 tations. First, the use of radiomic predictions can be hindered 183 by the cost of collecting imaging data (12). Furthermore, in-184 sufficient performance of the radiomic prediction can result 185 in more modest improvements in (or, in extreme cases, even 186 loss of) statistical power. Cost-benefit analyses are thus war-187 ranted. Finally, if a radiomic predictor is trained on data 188 from a different population compared with those studied in 189 the current trial, the improvements in statistical power may 190 be less pronounced. However, due to the randomization in 191 the study, the type I error rate is expected to be maintained 192 and internal validation or calibration of the predictive model 193 is possible using data from the control arm of a clinical trial. 194

Further studies of the misspecification of the predictive 195 model as well as the clinical trial outcome model are warranted 196 for assessing potential gains and loss of power in these settings. 197 However, misspecification of clinical outcome models can be 198 guarded against using statistical models satisfying symmetry 199 criteria (13). The biomarkers from the two cases presented 200 here were both built using SVMs, but this methodology can 201 accommodate predictions more generally. Incorporation of 202 these biomarkers into a one-arm trial designs in which all 203 participants in the trial are treated similarly also requires 204 further statistical research. 205

**ACKNOWLEDGMENTS.** This work was supported in part 206 by NIH grants R01NS085211, R01MH112847, R01NS060910, 207 RF1AG054409, R01EB022573, R01NS042645, and R01MH112070. 208

- RC Petersen, et al., Alzheimer's disease neuroimaging initiative (adni): clinical characterization. Neurology 74, 201–209 (2010).
- BH Menze, et al., The multimodal brain tumor image segmentation benchmark (brats). *IEEE transactions on medical imaging* 34, 1993–2024 (2014).

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

- C Davatzikos, F Xu, Y An, Y Fan, SM Resnick, Longitudinal progression of alzheimer's-like patterns of atrophy in normal older adults: the spare-ad index. *Brain* 132, 2026–2035 (2009).
- L Macyszyn, et al., Imaging patterns predict patient survival and molecular subtype in glioblastoma via machine learning techniques. *Neuro-oncology* 18, 417–425 (2015).
- C Davatzikos, P Bhatt, LM Shaw, KN Batmanghelich, JQ Trojanowski, Prediction of mci to ad conversion, via mri, csf biomarkers, and pattern classification. *Neurobiol. aging* 32, 2322–e19 (2011).
- SJ Pocock, The combination of randomized and historical controls in clinical trials. J. chronic diseases 29, 175–188 (1976).
- K Viele, et al., Use of historical control data for assessing treatment effects in clinical trials. *Pharm. statistics* 13, 41–54 (2014).
- A Abadie, A Diamond, J Hainmueller, Synthetic control methods for comparative case studies: Estimating the effect of california's tobacco control program. J. Am. statistical Assoc. 105, 493–505 (2010).
- X Da, et al., Integration and relative value of biomarkers for prediction of mci to ad progression: spatial patterns of brain atrophy, cognitive scores, apoe genotype and csf biomarkers. *NeuroImage: Clin.* 4, 164–173 (2014).
- Y Fan, et al., Spatial patterns of brain atrophy in mci patients, identified via high-dimensional pattern classification, predict subsequent cognitive decline. *Neuroimage* **39**, 1731–1743 (2008).
- PK Crane, et al., Development and assessment of a composite score for memory in the alzheimer's disease neuroimaging initiative (adni). Brain imaging behavior 6, 502–516 (2012).
- K Kantarci, CR Jack, Neuroimaging in alzheimer disease: an evidence-based review. Neuroimaging Clin. 13, 197–209 (2003).
- RT Shinohara, CE Frangakis, CG Lyketsos, A broad symmetry criterion for nonparametric validity of parametrically based tests in randomized trials. *Biometrics* 68, 85–91 (2012).