

The Prostate Cancer DREAM Challenge: A Community-Wide Effort to Use Open Clinical Trial Data for the Quantitative Prediction of Outcomes in Metastatic Prostate Cancer

KALD ABDALLAH,^a CHARLES HUGH-JONES,^b THEA NORMAN,^c STEPHEN FRIEND,^c GUSTAVO STOLOVITZKY^{d,e}

^aProject Data Sphere, Raleigh, North Carolina, USA; ^bSanofi US, Bridgewater, New Jersey, USA; ^cSage Bionetworks, Seattle, Washington, USA; ^dIBM T.J. Watson Research Center, Yorktown Heights, New York, USA; ^eIcahn School of Medicine at Mount Sinai, New York, New York, USA

Disclosures of potential conflicts of interest may be found at the end of this article.

Massive quantities of digital health care data have been generated over the past 10 years and are the perfect catalyst to foster an open ecosystem in medical research. With cloud computing infrastructure that can provide a homogeneous platform for distributing and processing data, it is now possible to draw on vast, diverse, and worldwide pools of human talent to generate insights for understanding disease and, ultimately, improving patient outcomes [1].

The use of open competitions—“crowdsourced” challenges—is an increasingly effective way to engage these diverse communities. Contributors come from multiple disciplines including statistics, machine learning, and computational biology. A number of initiatives have successfully run crowdsourced challenges in medical and basic research in which teams competitively solved specific problems within fixed time periods [2]. Notable examples include the Critical Assessment of Protein Structure Prediction (CASP) [3], the CLARITY Challenge for standardizing clinical genome sequencing analysis and reporting [4], and the DREAM Challenges (Dialogue for Reverse Engineering Assessments and Methods; <http://dreamchallenges.org>) [5] for the assessment of predictive models of disease.

Representative of the prevalence and potential of the model to accelerate research, the DREAM community has launched 32 challenges in the past 8 years, published >100 journal articles, and seen the participation of >8,000 solvers.

As a new paradigm for scientific research, crowdsourced challenges complement traditional methods. They have four key characteristics. First, they bring rigor to the field of predictive modeling in that final predictions are made from blind data sets, and this prevents common flaws related to such methodology (e.g., data leakage and overfitting). Second, the structure of a challenge is an optimal approach for comparing methodologies on the same basis and provides a guide for methods that efficiently extract relevant patterns from data. Third, if organized in an open and collaborative manner, challenges create communities of researchers that would not interact otherwise and increase access to data sets

that might otherwise be available to only a single laboratory. Fourth, the results of crowdsourced challenges, including algorithms, source code, and analysis, represent an important new resource for future research.

Despite the enormous potential of crowdsourced challenges to accelerate research, the hosting of clinical trial data for open challenges has been impeded by a lack of broad access to the quantities of clinical trial data generated, by difficulties in effective patient data deidentification, by intellectual property concerns, and by technological infrastructure costs.

Three years ago, DREAM and the nonprofit group Prize4Life addressed many of these issues [6]. They ran the Amyotrophic Lateral Sclerosis (ALS) Progression Challenge, a first-of-its-kind crowdsourced challenge using clinical trial data. The challenge asked participants to develop models to predict the speed of the progression of motor dysfunction in ALS patients, and leveraged data from several ALS clinical trials that had been consolidated in the Pooled Resource Open-Access ALS Clinical Trials database (PRO-ACT; <https://nctu.partners.org/ProACT>). The challenge drew 1,073 registrants from >60 countries. Importantly, the majority of challenge teams that participated had not previously worked on ALS.

The challenge resulted in the submission of 37 unique algorithms from which the top-performing models were identified. The improvement in predictability of disease progression translated to an estimated 20% decrease in the number of patients needed for enrollment in ALS trials. Furthermore, the winning algorithms could reliably discriminate patients with “fast” versus “slow” disease progression with better accuracy than expert clinicians who received the same data [6].

The ALS challenge demonstrated that crowdsourcing clinical data can enhance clinical research and potentially provide effective clinical measures of aggressiveness. In 2012, two further developments took place.

First, DREAM partnered its challenges with Sage Bionetworks (Seattle, WA, <http://sagebase.org>), a nonprofit company focused on developing open systems, incentives, and standards to accelerate data-driven predictive modeling for

biomedical research. Since 2012, DREAM challenges have been hosted on Synapse, Sage's institutional review board-approved open computational platform that facilitates sharing of sensitive data and that provides a working environment, registration services, and access to data, leaderboards, forums, and provenance-based data analysis.

Second, in 2012, the Life Sciences Consortium of the CEO Roundtable on Cancer (Cary, NC, <http://ceo-lsc.org>), another not-for profit group, conceived a data-sharing platform called Project Data Sphere (<https://www.projectdatasphere.org>). This platform launched in 2014 and is committed to broad sharing and analysis of historical comparator-arm phase III oncology data sets. This unprecedented effort brought together many organizations and corporations that collaboratively and systematically addressed historical barriers to broadly sharing cancer clinical trial data. Initial data sets in multiple tumor types were provided by six major pharmaceutical companies and by Memorial Sloan Kettering Cancer Center. These data are being used to fuel the Prostate Cancer DREAM Challenge.

The collaboration of Project Data Share and Sage Bionetworks/DREAM, in launching the Prostate Cancer DREAM Challenge, set the goal of improving a predictive model of disease progression and treatment toxicity in prostate cancer using historical trial data. Predictions identified through this challenge have the potential to translate into reduced trial redundancy, better clinical decision tools, and improved patient outcomes.

The process of developing the challenge was complex but reproducible. It required the collaboration of diverse contributors. A voluntary group of prostate cancer researchers assessed the challenge structure. Four data sets consisting of

docetaxel-treated phase III control arms were curated by data scientists. A rigorous evaluation ("dry run") validated the challenge questions, methodology, and process to objectively score contributions to identify the top-performing models. The challenge launches on March 16, 2015.

Although the potential for crowdsourcing in biomedical research is just beginning to be appreciated, its potential is considerable. An exponential increase in the volume and complexity of clinical data drives the necessity of collaboration. Large networks of diverse collaborators and technology and broad access to data are rapidly becoming the sine qua non for major scientific advances.

The Prostate Cancer DREAM Challenge includes all of these essential elements. It can serve as an important model for 21st century open clinical research while addressing the unmet needs of metastatic prostate cancer patients, whose overall mortality remains high.

AUTHOR CONTRIBUTIONS

Conception/Design: Kald Abdallah, Thea Norman, Gustavo Stolovitzky

Manuscript writing: Kald Abdallah, Charles Hugh-Jones, Thea Norman, Stephen Friend, Gustavo Stolovitzky

Final approval of manuscript: Kald Abdallah, Charles Hugh-Jones, Thea Norman, Stephen Friend, Gustavo Stolovitzky

DISCLOSURES

Kald Abdallah: Project Data Sphere LLC (E); **Charles Hugh-Jones:** Project Data Sphere LLC (C/A), Sanofi Aventis (E, OI). The other authors indicated no financial relationships.

(C/A) Consulting/advisory relationship; (RF) Research funding; (E) Employment; (ET) Expert testimony; (H) Honoraria received; (OI) Ownership interests; (IP) Intellectual property rights/inventor/patent holder; (SAB) Scientific advisory board

REFERENCES

1. Surowiecki J. *The Wisdom of Crowds*. New York, NY: Anchor Books, 2005.
2. Boutros PC, Margolin AA, Stuart JM et al. Toward better benchmarking: Challenge-based methods assessment in cancer genomics. *Genome Biol* 2014;15:462.
3. Shi S, Pei J, Sadreyev RI et al. Analysis of CASP8 targets, predictions and assessment methods. *Database (Oxford)* 2009;bap003.
4. Brownstein CA, Beggs AH, Homer N et al. An international effort towards developing standards for best practices in analysis, interpretation and reporting of clinical genome sequencing results in the CLARITY Challenge. *Genome Biol* 2014;15:R53.
5. Jarchum I, Jones S. DREAMing of benchmarks. *Nat Biotechnol* 2015;33:49–50.
6. Küffner R, Zach N, Norel R et al. Crowdsourced analysis of clinical trial data to predict amyotrophic lateral sclerosis progression. *Nat Biotechnol* 2015; 33:51–57.

EDITOR'S NOTE: See also in this issue two articles related to Project Data Sphere and its potential to advance the treatment of prostate cancer.

Angela K. Green, Katherine E. Reeder-Hayes, Robert W. Corty et al. The Project Data Sphere Initiative: Accelerating Cancer Research by Sharing Data (p. 464)

Angela K. Green, Robert W. Corty, William A. Wood et al. Comparative Effectiveness of Mitoxantrone Plus Prednisone Versus Prednisone Alone in Metastatic Castrate-Resistant Prostate Cancer After Docetaxel Failure (p. 516)