Optimizers manuscript

This manuscript (<u>permalink</u>) was automatically generated from <u>greenelab/optimizer-manuscript@b2b0834</u> on April 25, 2023.

Authors

Jake Crawford

(D) 0000-0001-6207-0782 · **(Ω)** jjc2718 · **У** jjc2718

Genomics and Computational Biology Graduate Group, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA

• Casey S. Greene [™]

Department of Biochemistry and Molecular Genetics, University of Colorado School of Medicine, Aurora, CO, USA; Center for Health AI, University of Colorado School of Medicine, Aurora, CO, USA

■ — Correspondence possible via <u>GitHub Issues</u> or email to Casey S. Greene <casey.s.greene@cuanschutz.edu>.

Abstract

Introduction

Gene expression profiles are widely used to classify samples or patients into relevant groups or categories, both preclinically [1,2] and clinically [3,4]. To extract informative gene features and to perform classification, a diverse array of algorithms exist, and different algorithms perform well across varying datasets and tasks [1]. Even within a given model class, multiple optimization methods can often be applied to find well-performing model parameters or to optimize a model's loss function. One commonly used example is logistic regression. The widely used scikit-learn Python package for machine learning [5] provides two modules for fitting logistic regression classifiers:

LogisticRegression, which uses the liblinear coordinate descent method [6] to find parameters that optimize the logistic loss function, and SGDClassifier, which uses stochastic gradient descent [7] to optimize the same loss function.

Using scikit-learn, we compared the liblinear (coordinate descent) and SGD optimization techniques for prediction of driver mutation status in tumor samples, across a wide variety of genes implicated in cancer initiation and development [8]. We applied LASSO (L1-regularized) logistic regression, and tuned the strength of the regularization to compare model selection between optimizers. We found that across a variety of models (i.e. varying regularization strengths), the training dynamics of the optimizers were considerably different: models fit using liblinear tended to perform best at fairly high regularization strengths (100-1000 nonzero features in the model) and overfit easily with low regularization strengths. On the other hand, models fit using stochastic gradient descent tended to perform best at fairly low regularization strengths (10000+ nonzero features in the model), and overfitting was uncommon.

Our results caution against viewing optimizer choice as a "black box" component of machine learning modeling. The observation that LASSO logistic regression models fit using SGD tended to perform best for low levels of regularization, across diverse driver genes, runs counter to conventional wisdom in statistics and machine learning for high-dimensional data which generally states that explicit regularization and/or feature selection is necessary. Comparing optimizers/model implementations directly is rare in applications of machine learning for genomics, and our work shows that this choice can affect generalization and interpretation properties of the model significantly. Based on our results, we recommend considering the appropriate optimization approach carefully based on the goals of each individual analysis.

Methods

Results

Discussion

References

1. The ability to classify patients based on gene-expression data varies by algorithm and performance metric

Stephen R Piccolo, Avery Mecham, Nathan P Golightly, Jérémie L Johnson, Dustin B Miller *PLOS Computational Biology* (2022-03-11) https://doi.org/gr43qd

DOI: <u>10.1371/journal.pcbi.1009926</u> · PMID: <u>35275931</u> · PMCID: <u>PMC8942277</u>

2. Supervised learning is an accurate method for network-based gene classification

Renming Liu, Christopher A Mancuso, Anna Yannakopoulos, Kayla A Johnson, Arjun Krishnan *Bioinformatics* (2020-04-14) https://doi.org/gmvnfc

DOI: 10.1093/bioinformatics/btaa150 · PMID: 32129827 · PMCID: PMC7267831

3. Supervised Risk Predictor of Breast Cancer Based on Intrinsic Subtypes

Joel S Parker, Michael Mullins, Maggie CU Cheang, Samuel Leung, David Voduc, Tammi Vickery, Sherri Davies, Christiane Fauron, Xiaping He, Zhiyuan Hu, ... Philip S Bernard *Journal of Clinical Oncology* (2009-03-10) https://doi.org/c2688w

DOI: <u>10.1200/jco.2008.18.1370</u> · PMID: <u>19204204</u> · PMCID: <u>PMC2667820</u>

4. Prediction of adjuvant chemotherapy benefit in endocrine responsive, early breast cancer using multigene assays

Kathy S Albain, Soonmyung Paik, Laura van't Veer

The Breast (2009-10) https://doi.org/bp4rtw

DOI: <u>10.1016/s0960-9776(09)70290-5</u> · PMID: <u>19914534</u>

5. Scikit-learn: Machine Learning in Python

Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, ... Édouard Duchesnay

Journal of Machine Learning Research (2011) http://jmlr.org/papers/v12/pedregosa11a.html

6. LIBLINEAR: A Library for Large Linear Classification

Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, Chih-Jen Lin *Journal of Machine Learning Research* (2008) http://jmlr.org/papers/v9/fan08a.html

7. Online Learning and Stochastic Approximations

Leon Bottou

(1998) https://wiki.eecs.yorku.ca/course archive/2012-13/F/6328/ media/bottou-onlinelearning-98.pdf

8. Cancer Genome Landscapes

B Vogelstein, N Papadopoulos, VE Velculescu, S Zhou, LA Diaz, KW Kinzler *Science* (2013-03-28) https://doi.org/6rg

DOI: 10.1126/science.1235122 · PMID: 23539594 · PMCID: PMC3749880