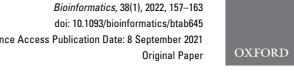
Downloaded from https://academic.oup.com/bioinformatics/article/38/1/157/6366546 by Columbia University user on 13 January 2022



Gene expression

Learning sparse log-ratios for high-throughput sequencing data

Elliott Gordon-Rodriguez (b) 1,*, Thomas P. Quinn (b) 2 and John P. Cunningham¹

¹Department of Statistics, Columbia University, New York, NY 10025, USA and ²Applied Artificial Intelligence Institute, Deakin University, Geelong, VIC 3126, Australia

*To whom correspondence should be addressed.

Associate Editor: Pier Luigi Martelli

Received on June 6, 2021; revised on August 9, 2021; editorial decision on September 1, 2021; accepted on September 3, 2021

Abstract

Motivation: The automatic discovery of sparse biomarkers that are associated with an outcome of interest is a central goal of bioinformatics. In the context of high-throughput sequencing (HTS) data, and compositional data (CoDa) more generally, an important class of biomarkers are the log-ratios between the input variables. However, identifying predictive log-ratio biomarkers from HTS data is a combinatorial optimization problem, which is computationally challenging. Existing methods are slow to run and scale poorly with the dimension of the input, which has limited their application to low- and moderate-dimensional metagenomic datasets.

Results: Building on recent advances from the field of deep learning, we present CoDaCoRe, a novel learning algorithm that identifies sparse, interpretable and predictive log-ratio biomarkers. Our algorithm exploits a continuous relaxation to approximate the underlying combinatorial optimization problem. This relaxation can then be optimized efficiently using the modern ML toolbox, in particular, gradient descent. As a result, CoDaCoRe runs several orders of magnitude faster than competing methods, all while achieving state-of-the-art performance in terms of predictive accuracy and sparsity. We verify the outperformance of CoDaCoRe across a wide range of microbiome, metabolite and microRNA benchmark datasets, as well as a particularly high-dimensional dataset that is outright computationally intractable for existing sparse log-ratio selection methods.

Availability and implementation: The CoDaCoRe package is available at https://github.com/egr95/R-codacore, Code and instructions for reproducing our results are available at https://github.com/cunningham-lab/codacore.

Contact: eg2912@columbia.edu

Supplementary information: Supplementary data are available at Bioinformatics online.

1 Introduction

High-throughput sequencing (HTS) technologies have enabled the relative quantification of the different bacteria, metabolites or genes, that are present in a biological sample. However, the nature of these recording technologies results in sequencing biases that complicate the analysis of HTS data. In particular, HTS data come as counts, whose totals are constrained to the capacity of the measuring device. These totals are an artifact of the measurement process, and do not depend on the subject being measured. Hence, HTS counts arguably should be interpreted in terms of relative abundance; in statistical terminology, it follows that HTS data are an instance of compositional data (CoDa) (Calle, 2019; Gloor et al., 2016, 2017; Quinn et al., 2018, 2019).

Mathematically, CoDa can be defined as a set of non-negative vectors whose totals are uninformative. Since the seminal work of Aitchison (1982), the statistical analysis of CoDa has become a discipline in its own right (Pawlowsky-Glahn and Buccianti, 2011; Pawlowsky-Glahn and Egozcue, 2006). But why does CoDa deserve special treatment? Unlike unconstrained real-valued data, the compositional nature of CoDa results in each variable becoming negatively correlated to all others (increasing one component of a composition implies a relative decrease of the other components). It is well known that, as a result, the usual measures of association and feature attribution are problematic when applied to CoDa (Filzmoser et al., 2009; Lovell et al., 2015; Pearson, 1896). Consequently, bespoke methods are necessary for a valid statistical analysis (Gloor et al., 2017). Indeed, the application of CoDa methodology to HTS data, especially microbiome data, has become increasingly popular in recent years (Calle, 2019; Fernandes et al., 2013, 2014; Quinn et al., 2021; Rivera-Pinto et al., 2018).

The standard approach for analyzing CoDa is based on applying log-ratio transformations to map our data onto unconstrained Euclidean space, where the usual tools of statistical learning apply (Pawlowsky-Glahn and Egozcue, 2006). The choice of the log-ratio transform offers the necessary property of scale invariance, but in 158 E.Gordon-Rodriguez et al.

the CoDa literature, it holds primacy for a variety of other technical reasons, including subcompositional coherence (Aitchison, 1982; Pawlowsky-Glahn and Buccianti, 2011). Log-ratios can be taken over pairs of input variables (Aitchison, 1982; Bates and Tibshirani, 2019; Greenacre, 2019b) or aggregations thereof, typically geometric means (Aitchison, 1982; Egozcue, 2003; Egozcue and Pawlowsky-Glahn, 2005; Rivera-Pinto et al., 2018) or summations (Greenacre, 2019a, 2020; Quinn and Erb, 2020). The resulting features work well empirically, but also imply a clear interpretation: a log-ratio is a single composite score that expresses the overall quantity of one sub-population as compared with another. For example, in microbiome HTS data, the relative weights between subpopulations of related microorganisms are commonly used as clinical biomarkers (Crovesy et al., 2020; Magne et al., 2020; Rahat-Rozenbloom et al., 2014). When the log-ratios are sparse, meaning they are taken over a small number of input variables, they define biomarkers that are particularly intuitive to understand, a key desiderata for predictive models that are of clinical relevance (Goodman and Flaxman, 2017).

Thus, learning sparse log-ratios is a central problem in CoDa. This problem is especially challenging in the context of HTS data, due to its high dimensionality (ranging from 100 to over 10 000 variables). Existing methods rely on stepwise search (Greenacre, 2019b; Rivera-Pinto et al., 2018) or evolutionary algorithms (Prifti et al., 2020; Quinn and Erb, 2020), which scale poorly with the dimension of the input. These algorithms are prohibitively slow for most HTS datasets, and thus there is a new demand for sparse and interpretable models that scale to high dimensions (Cammarota et al., 2020; Li, 2015; Susin et al., 2020)

This demand motivates the present work, in which we present CoDaCoRe, a novel learning algorithm for Compositional Data via Continuous Relaxations. CoDaCoRe builds on recent advances from the deep learning literature on *continuous relaxations* of discrete latent variables (Jang *et al.*, 2017; Linderman *et al.*, 2018); we design a novel relaxation that approximates a combinatorial optimization problem over the set of log-ratios. In turn, this approximation can be optimized efficiently using gradient descent, and subsequently discretized to produce a sparse log-ratio biomarker, thus dramatically reducing runtime without sacrificing interpretability nor predictive accuracy. The main contributions of our method can be summarized as follows:

- Computational efficiency. CoDaCoRe scales linearly with the dimension of the input. It runs several orders of magnitude faster than its competitors.
- Interpretability. CoDaCoRe identifies a set of log-ratios that are sparse, biologically meaningful and ranked in order of importance. Our model is highly interpretable, and much sparser, relative to competing methods of similar accuracy and computational complexity.
- Predictive accuracy. CoDaCoRe achieves better out-of-sample accuracy than existing CoDa methods, and performs similarly to state-of-the-art black-box classifiers (which are neither sparse nor interpretable).
- Ease of use. We devise an adaptive learning rate scheme that enables CoDaCoRe to converge reliably, requiring no additional hyperparameter tuning.

2 Background

Our work focuses on the supervised learning problem $\mathbf{x}_i \mapsto y_i$, where the inputs \mathbf{x}_i are HTS data (or any CoDa), and the outputs y_i are the outcome of interest. For many microbiome applications, \mathbf{x}_i represents a vector of frequencies of the different species of bacteria that compose the microbiome of the ith subject. In other words, x_{ij} denotes the abundance of the ith species (of which there are p total) in the ith subject. The response y_i is often a binary variable indicating whether the ith subject belongs to the case or the control groups

(e.g. sick versus healthy). For HTS data, the input frequencies x_{ij} arise from an inexhaustive sampling procedure, so that the totals $\sum_{j=1}^{p} x_{ij}$ are arbitrary and the components should only be interpreted in relative terms (i.e. as CoDa) (Calle, 2019; Gloor *et al.*, 2017; Gloor and Reid, 2016; Quinn *et al.*, 2018). While many of our applications pertain to microbiome data, our method applies to any high-dimensional HTS data, including those produced by *Liquid Chromatography Mass Spectrometry* (Filzmoser and Walczak, 2014).

2.1 Log-ratio analysis

Our goal is to obtain sparse log-ratio transformed features that can be passed to a downstream classifier or regression function. As discussed, these log-ratios will result in interpretable features and scaleinvariant models (that are also subcompositionally coherent), thus satisfying the key requirements for valid statistical inference in the context of CoDa. The simplest such choice is the pairwise log-ratio, defined as $\log(x_{ij^+}/x_{ij^-})$, where j^+ and j^- denote the indexes of a pair of input variables (Aitchison, 1982). Note that the ratio cancels out any scaling factor applied to x_i , preserving only the relative information, while the log transformation ensures the output is (unconstrained) real-valued. There are many such (j^+, j^-) pairs (to be precise, $p(p-1)/2 = O(p^2)$ of them). In order to select good pairwise log-ratios from a set of input variables, Greenacre (2019b) proposed a greedy step-wise search algorithm. This method produces a sparse and interpretable set of features, but it is prohibitively slow on high-dimensional datasets, as a result of the step-wise algorithm scaling quadratically in the dimension of the input. A heuristic search algorithm that is less accurate but computationally faster has been developed as part of Quinn et al. (2017), though its computational cost is still troublesome (as we shall see in Section 4). The logratio lasso is a computationally efficient alternative for selecting pairwise log-ratios (Bates and Tibshirani, 2019).

2.1.1 Balances

Recently, a class of log-ratios known as *balances* (Egozcue and Pawlowsky-Glahn, 2005) have become of interest in microbiome applications, due to their interpretability as the relative weight between two sub-populations of bacteria (Morton *et al.*, 2019; Quinn and Erb, 2019). Balances are defined as the log-ratios between geometric means of two subsets of the input variables (Note that the original definition of balances includes a "normalization" constant, which we omit for clarity. This constant is in fact unnecessary, as it will get absorbed into a regression coefficient downstream.):

$$B(\mathbf{x}_{i}; J^{+}, J^{-}) = \log \left(\frac{(\prod_{j \in J^{+}} x_{ij})^{\frac{1}{p^{+}}}}{(\prod_{j \in J^{-}} x_{ij})^{\frac{1}{p^{-}}}} \right)$$
(1)

$$= \frac{1}{p^+} \sum_{j \in J^+} \log x_{ij} - \frac{1}{p^-} \sum_{j \in J^-} \log x_{ij},$$

where J^+ and J^- denote a pair of disjoint subsets of the indices $\{1,\ldots,p\}$, and p^+ and p^- denote their respective sizes. For example, in microbiome data, J^+ and J^- are groups of bacteria species that may be related by their environmental niche (Morton *et al.*, 2017) or genetic similarity (Silverman *et al.*, 2017; Washburne *et al.*, 2017). Note that when $p^+ = p^- = 1$ (i.e. J^+ and J^- each contain a single element), $B(\mathbf{x}; J^+, J^-)$ reduces to a pairwise log-ratio. By allowing for the aggregation of more than one variable in the numerator and denominator of the log-ratio, balances provide a far richer set of features that allows for more flexible models than pairwise log-ratios. Insofar as the balances are taken over a small number of variables (i.e. J^+ and J^- are sparse), they also provide highly interpretable biomarkers.

The *selbal* algorithm (Rivera-Pinto *et al.*, 2018) has gained popularity as a method for automatically identifying balances that predict a response variable. However, this algorithm is also based on a

greedy step-wise search through the combinatorial space of subset pairs (J^+, J^-) , which scales poorly in the dimension of the input and becomes prohibitively slow for many HTS datasets (Susin *et al.*, 2020).

2.1.2 Amalgamations

An alternative to balances, known as *amalgamation*, is defined by aggregating components through summation:

$$A(\mathbf{x}_i; J^+, J^-) = \log\left(\frac{\sum_{j \in J^+} x_{ij}}{\sum_{j \in J^-} x_{ij}}\right),\tag{2}$$

where again J^+ and J^- denote disjoint subsets of the input components. Amalgamations have the advantage of reducing the dimensionality of the data through an operation, the sum, that some authors argue is more interpretable than a geometric mean (Greenacre, 2019a; Greenacre *et al.*, 2020). On the other hand, amalgamations can be less effective than balances for identifying components that are statistically important, but small in magnitude, e.g. rare bacteria species (since small terms will have less impact on a summation than on a product).

Recently, Greenacre (2020) has advocated for the use of expert-driven amalgamations, using domain knowledge to construct the relevant features. On the other hand, Quinn and Erb (2020) proposed *amalgam*, an evolutionary algorithm to automatically identify amalgamated log-ratios (Equation 2) that are predictive of a response variable. However, this algorithm does not scale to high-dimensional data (albeit, comparing favorably to selbal), nor does it produce sparse models (hindering interpretability of the results). A similar evolutionary algorithm can be found in Prifti *et al.* (2020), however, their model is not scale invariant, as is required by most authors in the field (Pawlowsky-Glahn and Egozcue, 2006).

Other relevant log-ratio methodology is briefly reviewed in Supplementary Material (Supplementary Section SB).

3 Materials and methods

We now present CoDaCoRe, a novel learning algorithm for HTS data, and more generally, high-dimensional CoDa. Unlike existing methods, CoDaCoRe is simultaneously scalable, interpretable, sparse and accurate. In Table 1 from Supplementary Material (Supplementary Section SC), we summarize the relative merits of CoDaCoRe and its competitors.

3.1 Optimization problem

In its basic formulation, CoDaCoRe learns a regression function of the form:

$$f(\mathbf{x}) = \alpha + \beta \cdot B(\mathbf{x}; J^+, J^-), \tag{3}$$

where B denotes a balance (Equation 1), and α and β are scalar parameters. This regression function can be thought of in two stages: (i) we take the input and use it to compute a balance score and (ii) we feed the balance score to a logistic regression classifier. For clarity, we will restrict our exposition to this formulation, but note that our algorithm can be applied equally to learn amalgamations instead of balances (see Section 3.6), as well as generalizing straightforwardly to non-linear functions (provided they are suitably parameterized and differentiable).

Let L(y,f) denote the cross-entropy loss, with $f \in \mathbb{R}$ given in logit space. The goal of CoDaCoRe is to find the balance that is maximally associated with the response. Mathematically, this can be written as:

$$\min_{(J^+J^-,\alpha,\beta)} \sum_{i} L\Big(y_i, \alpha + \beta \cdot B(\mathbf{x}_i; J^+, J^-)\Big). \tag{4}$$

This objective function may look similar to a univariate logistic regression, however, our problem is complicated by the joint optimization over the subsets J^+ and J^- which determine the input variables that compose the balance. Note that the number of possible subsets of p variables is 2^p , so the set of possible balances is greater than 2^p and grows *exponentially* in p. Exact optimization is therefore computationally intractable for any but the smallest of datasets, and an approximate solution is required. Selbal corresponds to one such approximation, offering *quadratic* complexity in p, which is practical for low- to moderate-dimensional datasets (p < 100), but does not scale to high dimensions (p > 1000). As we shall now see, CoDaCoRe represents a critical improvement, achieving *linear* complexity in p which dramatically reduces runtime and enables, for the first time, the use of balances and amalgamations for the analysis of high-dimensional HTS data.

3.2 Continuous relaxation

The key insight of CoDaCoRe is to approximate our combinatorial optimization problem (Equation 4) with a continuous relaxation that can be trained efficiently by gradient descent. Our relaxation is inspired by recent advances in deep learning models with discrete latent variables (Jang et al., 2017; Linderman et al., 2018; Maddison et al., 2017; Mena et al., 2018; Potapczynski et al., 2020). However, we are not aware of any similar proposals for optimizing over

 $\textbf{Table 1}. \ \textbf{Evaluation metrics shown for each method, averaged over 25 \ datasets} \times \textbf{20} \ \textbf{random train/test splits}$

	Runtime (s)	Active inputs (%)	Accuracy (%)	AUC (%)	F1 (%)
CoDaCoRe—Balances (ours)	4.5 ±0.4	1.9 ±0.3	75.2 ± 2.4	79.5 ± 2.6	73.7 ± 2.6
CoDaCoRe—Amalgamations (ours)	4.4 ± 0.4	1.9 ± 0.3	71.8 ± 2.4	74.5 ± 2.8	69.8 ± 2.9
selbal (Rivera-Pinto et al., 2018)	$79\ 033.7 \pm 2094.1$	2.4 ± 0.2	61.2 ± 1.9	80.0 ± 2.4	70.9 ± 1.1
Pairwise log-ratios (Greenacre, 2019b)	$14\ 207.0 \pm 1038.4$	2.5 ± 0.4	73.3 ± 1.7	75.2 ± 2.4	67.8 ± 3.0
Lasso	1.6 ± 0.1	4.4 ± 0.6	72.4 ± 1.7	75.2 ± 2.3	65.2 ± 3.7
CoDaCoRe—balances with $\lambda = 0$ (ours)	9.8 ± 2.2	6.1 ± 0.7	77.6 ± 2.2	82.0 ± 2.3	76.0 ± 2.5
Coda-lasso (Lu et al., 2019)	1043.0 ± 55.4	19.7 ± 2.7	72.5 ± 2.3	78.0 ± 2.4	64.2 ± 4.4
amalgam (Quinn and Erb, 2020)	7360.5 ± 209.8	87.6 ± 2.1	74.4 ± 2.5	78.2 ± 2.7	73.9 ± 2.8
DeepCoDA (Quinn et al., 2020)	296.5 ± 21.4	89.3 ± 0.6	70.6 ± 2.9	77.6 ± 2.9	64.7 ± 7.4
CLR-lasso (Susin et al., 2020)	2.0 ± 0.2	100.0 ± 0.0	77.5 ± 1.8	81.6 ± 2.2	75.8 ± 2.7
Random Forest	10.6 ± 0.4	_	78.0 ± 2.2	82.2 ± 2.2	77.3 ± 2.5
Log-ratio lasso (Bates and Tibshirani, 2019)*	135.0 ± 11.1	0.7 ± 0.0	72.0 ± 2.4	76.4 ± 2.3	69.2 ± 2.7

Note: Standard errors are computed independently on each dataset, and then averaged over the 25 datasets. The models are ordered by sparsity, i.e. percentage of active input variables. CoDaCoRe (with balances) is the only learning algorithm that is simultaneously fast, sparse and accurate. The penultimate row shows the performance of Random Forest, a powerful black-box classifier which can be thought of as providing an approximate upper bound on the predictive accuracy of any interpretable model. The bottom row is shown separately and marked with an asterisk because the corresponding algorithm failed to converge on 432 out our 500 runs (averages were taken after imputing these missing values with the corresponding values obtained with pairwise log-ratios, which is the most similar method). We highlight in bold the CoDa models that are fast to run, as well as the CoDa models that are most sparse and accurate.

160 E.Gordon-Rodriguez et al.

disjoint subsets, nor for learning balances or amalgamations in the context of CoDa.

Our relaxation is parameterized by an unconstrained vector of 'assignment weights', $\mathbf{w} \in \mathbb{R}^p$, with one scalar parameter per input dimension (e.g. one weight per bacteria species). The weights are mapped to a vector of 'soft assignments' via:

$$\tilde{\mathbf{w}} = 2 \cdot \operatorname{sigmoid}(\mathbf{w}) - 1 = \frac{2}{1 + \exp(-\mathbf{w})} - 1, \tag{5}$$

where the sigmoid is applied component-wise. Intuitively, large positive weights will max out the sigmoid, leading to soft assignments close to +1, whereas large negative weights will zero out the sigmoid, resulting in soft assignments close to -1. This mapping is akin to softly assigning input variables to the groups J^+ and J^- respectively.

Let us write $\tilde{\mathbf{w}}^+ = \text{ReLU}(\tilde{\mathbf{w}})$ and $\tilde{\mathbf{w}}^- = \text{ReLU}(-\tilde{\mathbf{w}})$ for the (component-wise) positive and negative parts of $\tilde{\mathbf{w}}$, respectively. We approximate balances (Equation 1) with the following relaxation:

$$\tilde{B}(\mathbf{x}_i; \mathbf{w}) = \frac{\sum_j \tilde{w}_j^+ \log x_{ij}}{\sum_j \tilde{w}_j^+} - \frac{\sum_j \tilde{w}_j^- \log x_{ij}}{\sum_j \tilde{w}_j^-}$$
(6)

$$= \frac{\tilde{\mathbf{w}}^+ \cdot \log \mathbf{x}_i}{||\tilde{\mathbf{w}}^+||_1} - \frac{\tilde{\mathbf{w}}^- \cdot \log \mathbf{x}_i}{||\tilde{\mathbf{w}}^-||_1}. \tag{7}$$

In other words, we approximate the geometric averages over subsets of the inputs, by *weighted* geometric averages over all components (compare Equations 1 and 6).

Crucially, this relaxation is differentiable in \mathbf{w} , allowing us to construct a surrogate objective function that can be optimized jointly in $(\mathbf{w}, \alpha, \beta)$ by gradient descent:

$$\min_{(\mathbf{w},\alpha,\beta)} \sum_{i} L(\mathbf{y}_{i}, \alpha + \beta \cdot \tilde{B}(\mathbf{x}_{i}; \mathbf{w})). \tag{8}$$

Moreover, the computational cost of differentiating this objective function scales linearly in the dimension of w, which overall results in linear scaling for our algorithm. We also note that the functional form of our relaxation (Equation 6) can be exploited in order to select the learning rate adaptively (i.e. without tuning), resulting in robust convergence across all real and simulated datasets that we considered. We defer the details of our implementation of gradient descent to Supplementary Material (Supplementary Section SC.1).

3.3 Discretization

While a set of features in the form of Equation 6 may perform accurate classification, a weighted geometric average over all input variables is much harder for a biologist to interpret (and less intuitively appealing) than a bona fide balance over a small number of variables. For this reason, CoDaCoRe implements a 'discretization' procedure that exploits the information learned by the soft assignment vector $\tilde{\mathbf{w}}$, in order to efficiently identify a pair of sparse subsets, \hat{J}^+ and \hat{J}^- , which will define a balance.

The most straightforward way to convert the (soft) assignment $\tilde{\mathbf{w}}$ into a (hard) pair of subsets is by fixing a threshold $t \in (0, 1)$:

$$\tilde{J}^+ = \{ j : \tilde{w}_i > t \}, \tag{9}$$

$$\tilde{J}^- = \{j : \tilde{w}_j < -t\}.$$
 (10)

Note that given a trained $\tilde{\mathbf{w}}$ and a fixed threshold t, we can evaluate the quality of the corresponding balance $B(\mathbf{x}; \tilde{f}^+, \tilde{f}^-)$ (resp. amalgamation) by optimizing Equation 4 over (α, β) alone, i.e. fitting a linear model. Computationally, fitting a linear model is much faster than optimizing Equation 8, and can be done repeatedly for a range of values of t with little overhead. In CoDaCoRe, we combine this strategy with cross-validation in order to select the threshold, \hat{t} , that optimizes predictive performance (see Supplementary Section

SC.2 of Supplementary Material for full detail). Finally, the trained regression function is:

$$\hat{f}(\mathbf{x}) = \hat{\alpha} + \hat{\beta} \cdot B(\mathbf{x}; \hat{J}^+, \hat{J}^-), \tag{11}$$

where \hat{f}^+ and \hat{f}^- are the subsets corresponding to the optimal threshold \hat{t} , and $(\hat{\alpha}, \hat{\beta})$ are the coefficients obtained by regressing y_i against $B(\mathbf{x}_i; \hat{f}^+, \hat{f}^-)$ on the entire training set.

3.4 Regularization

Note from Equations 9 and 10 that larger values of t result in fewer input variables assigned to the balance $B(\mathbf{x}; \vec{J}^+, \vec{J}^-)$, i.e. a sparser model. Thus, CoDaCoRe can be regularized simply by making \hat{t} larger. Similar to lasso regression, CoDaCoRe uses the 1-standard*error* rule: namely, to pick the sparsest model (i.e. the highest *t*) with mean cross-validated score within 1 standard error of the optimum (Friedman et al., 2001). Trivially, this rule can be generalized to a λ standard-error rule (to pick the sparsest model within λ standard errors of the optimum), where λ becomes a regularization hyperparameter that can be tuned by the practitioner if so desired (with lower values trading off some sparsity in exchange for predictive accuracy). In our public implementation, $\lambda = 1$ is our default value, and this is used throughout our experiments (except where we indicate otherwise). In practice, lower values (e.g. $\lambda = 0$) can be useful when the emphasis is on predictive accuracy rather than interpretability or sparsity, though our benchmarks showed competitive performance for any $\lambda \in [0, 1]$.

3.5 CoDaCoRe algorithm

The computational efficiency of our continuous relaxation allows us to train multiple regressors of the form of Equation 11 within a single model. In the full CoDaCoRe algorithm, we ensemble multiple such regressors in a stage-wise additive fashion, where each successive balance is fitted on the residual from the current model. Thus, CoDaCoRe identifies a sequence of balances, in decreasing order of importance, each of which is sparse and interpretable. Training terminates when an additional relaxation (Equation 6) cannot improve the cross-validation score relative to the existing ensemble (equivalently, when we obtain $\hat{t} = 1$). Typically, only a small number of balances is required to capture the signal in the data, and as a result CoDaCoRe produces very sparse models overall, further enhancing interpretability. In Supplementary Material, we summarize our procedure in Supplementary Algorithm S1 (Supplementary Section SC) and we describe a number of extensions to the CoDaCoRe framework (Supplementary Section SD), including unsupervised learning.

3.6 Amalgamations

CoDaCoRe can be used to learn amalgamations (Equation 2) much in the same way as for balances (the choice of which to use depending on the goals of the biologist). In this case, our relaxation is defined as:

$$\tilde{A}(\mathbf{x}_i; \mathbf{w}) = \log \left(\frac{\sum_j \tilde{w}_j^+ x_{ij}}{\sum_j \tilde{w}_j^- x_{ij}} \right) = \log \left(\frac{\tilde{\mathbf{w}}^+ \cdot \mathbf{x}_i}{\tilde{\mathbf{w}}^- \cdot \mathbf{x}_i} \right), \tag{12}$$

i.e. we approximate summations over subsets of the inputs, with weighted summations over all components (compare Equation 2 and Equation 12). The rest of the argument follows verbatim, replacing $B(\cdot)$ with $A(\cdot)$ and $\tilde{B}(\cdot)$ with $\tilde{A}(\cdot)$ in Equations 3, 4, 8 and 11.

4 Experiments

We evaluate CoDaCoRe on a collection of 25 benchmark datasets including 13 datasets from the *Microbiome Learning Repo* (Vangay *et al.*, 2019), and 12 microbiome, metabolite and microRNA datasets curated by Quinn and Erb (2019). These data vary in dimension from 48 to 3090 input variables (see Supplementary Section SE of Supplementary Material for a full description). For each dataset, we fit CoDaCoRe and competing methods on 20 random 80/20 train/

test splits, sampled with stratification by case–control (He and Ma, 2013). Competing methods and their implementation are described in Supplementary Section SF.2 of Supplementary Material.

4.1 Results

We evaluate the quality of our models across the following criteria: computational efficiency (as measured by runtime), sparsity (as measured by the percentage of input variables that are active in the model) and predictive accuracy (as measured by out-of-sample accuracy, ROC AUC and F1 score). Table 1 provides an aggregated summary of the results; CoDaCoRe (with balances) is performant on all metrics. Indeed, our method provides the only interpretable model that is simultaneously scalable, sparse and accurate. Detailed performance metrics on each of the 25 datasets are provided in Supplementary Section SF of Supplementary Material, together with critical difference diagrams for each of our success metrics.

Figure 1 shows the average runtime of our classifiers on each dataset, with larger points denoting larger datasets. On these common benchmark datasets, CoDaCoRe trains up to 5 orders of magnitude faster than existing interpretable CoDa methods. On our larger datasets (3090 inputs), selbal runs in ~100 hours, pairwise log-ratios and amalgam both run in ~10 hours, and CoDaCoRe runs in under 10 seconds (full runtimes are provided in Supplementary Table S5 in Supplementary Material). All runs, including those involving gradient descent, were performed on identical CPU cores; CoDaCoRe can be accelerated further using GPUs, but we did not find it necessary to do so. It is also worth noting that the outperformance of CoDaCoRe is not merely as a result of the other methods failing on high-dimensional datasets. Supplementary Section SF.1.1 and Supplementary Figure S4 in Supplementary Material show that CoDaCoRe performs consistently across low- and high-dimensional datasets, and enjoys better sample efficiency than competing methods. Better sample efficiency could represent a particular advantage in biomedical studies, where most datasets have low n and high p.

Not only is CoDaCoRe sparser and more accurate than other interpretable models, it also performs on par with state-of-the-art black-box classifiers. By simply reducing the regularization parameter, from $\lambda=1$ to $\lambda=0$, CoDaCoRe (with balances) achieved an average 77.6% out-of-sample accuracy of and 82.0% AUC, on par with Random Forest (penultimate row of Table 1), while only using 5.9% of the input variables, on average. This result indicates, first, that CoDaCoRe provides a highly effective algorithm for variable

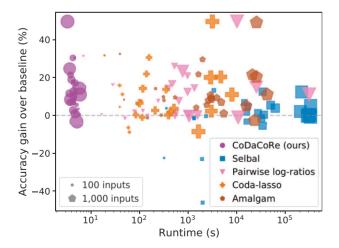


Fig. 1. Gain in classification accuracy (relative to the "majority vote" baseline classifier) plotted against runtime. Each point represents one of 25 datasets, with size proportional to the input dimension. Note the x-axis is drawn on the log-scale. CoDaCoRe (with balances) is the only method that scales effectively to our larger datasets, while consistently achieving high predictive accuracy. Moreover, its performance is broadly consistent across smaller and larger datasets

Table 2. Evaluation metrics for the liquid biopsy data (Best *et al.*, 2015), averaged over 20 independent 80/20 train/test splits

		Runtime (s)	Vars (#)	Acc. (%)	AUC (%)	F1 (%)
	CoDaCoRe	31±2.2	3 ± 1	91.0 ± 1.9	93.6 ± 2.6	94.4 ± 1.2
	Lasso	23 ± 0.2	22 ± 4	87.8 ± 1.3	94.7 ± 1.5	92.7 ± 0.7
	RF	383 ± 8.6	-	89.0 ± 1.6	94.1 ± 1.8	93.1 ± 1.0
	XGBoost	108 ± 1.6	-	90.6 ± 1.9	95.9 ± 1.5	94.1 ± 1.1

Note: CoDaCoRe (with balances) achieves equal predictive accuracy as competing methods, but with much sparser solutions. Note that sparsity is expressed as an (integer) number of active variables in the model (not as a percentage of the total, as was done in Table 1). We highlight in bold the sparsest and most accurate models.

selection in high-dimensional HTS data. Second, the fact that CoDaCoRe achieves similar predictive accuracy as state-of-the-art black-box classifiers, suggests that our model may have captured a near-complete representation of the signal in the data. At any rate, we take this as evidence that log-ratio transformed features are indeed of biological importance in the context of HTS data, corroborating previous microbiome research (Crovesy *et al.*, 2020; Magne *et al.*, 2020; Rahat-Rozenbloom *et al.*, 2014).

4.2 Interpretability

The CoDaCoRe algorithm offers two kinds of interpretability. First, it provides the analyst with sets of input variables whose aggregated ratio predicts the outcome of interest. These sets are easy to understand because they are discrete, with each component making an equivalent (unweighted) contribution. They are also sparse, usually containing fewer than 10 features per ratio, and can be made sparser by adjusting the regularization parameter λ. Such ratios have a precedent in microbiome research, for example the Firmicutes-to-Bacteroidetes ratio is used as a biomarker of gut health (Crovesy et al., 2020; Magne et al., 2020). Second, CoDaCoRe ranks predictive ratios hierarchically. Due to the ensembling procedure, the first ratio learned is the most predictive, the second ratio predicts the residual from the first, and so forth. Like principal components, the balances (or amalgamations) learned by CoDaCoRe are naturally ordered in terms of their explanatory power. This ordering aids

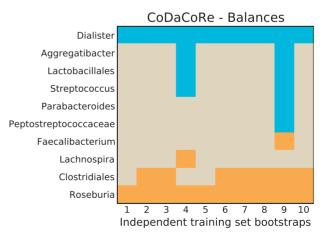


Fig. 2. CoDaCoRe variable selection for the first (most explanatory) log-ratio on the Crohn disease data (Rivera-Pinto et al., 2018). For each of 10 independent bootstraps of the training set (80% of the data randomly sampled with stratification by case—control), we show which variables are selected in the numerator (blue) and denominator (orange) of the balance. CoDaCoRe learns remarkably consistent logratios across independent training sets

162 E.Gordon-Rodriguez et al.

interpretability by decomposing a multivariable model into comprehensible 'chunks' of information.

Notably, we find a high degree of stability in the log-ratios selected by the model. We repeated CoDaCoRe on 10 independent training set splits of the Crohn disease data provided by Rivera-Pinto et al. (2018), and found consensus among the learned models. Figure 2 shows which bacteria were included for each split. Importantly, the bacteria that were selected consistently by CoDaCoRe—notably Dialister, Roseburia and Clostridiales—were also identified by Rivera-Pinto et al. (2018). In Supplementary Material (Supplementary Section SF), we also present a comparison of Figure 2 when using CoDaCoRe to learn amalgamations instead of balances. The amalgamations tend to select more abundant bacteria species like Faecalibacterium rather than rarer species like Roseburia (due to the geometric mean being more sensitive to small numbers than the summation operator).

4.3 Scaling to liquid biopsy data

HTS data generated from from clinical blood samples can be described as a 'liquid biopsy' that can be used for cancer diagnosis and surveillance (Alix-Panabières and Pantel, 2016; Best et al., 2015). These data can be very high-dimensional, especially when they include all gene transcripts as input variables. In a clinical context, the use of log-ratio predictors is an attractive option because they automatically correct for inter-sample sequencing biases that might otherwise limit the generalizability of the models (Dillies et al., 2013). Unfortunately, existing log-ratio methods like selbal and amalgam simply cannot scale to liquid biopsy datasets that contain as many as 50 000 or more input variables.

The large dimensionality of such data has restricted its analysis to overly simplistic linear models, black-box models that are scalable but not interpretable, or suboptimal hybrid approaches where input variables must be pre-selected based on univariate measures (Best et al., 2015; Sheng et al., 2018; Zhang et al., 2017). Owing to its linear scaling, CoDaCoRe can be fitted to these data at a similar computational cost to a single lasso regression, i.e. under a minute on a single CPU core. Thus, CoDaCoRe can be used to discover interpretable and predictive log-ratios that are suitable for liquid biopsy cancer diagnostics, among other similar applications.

We showcase the capabilities of CoDaCoRe in this high-dimensional setting, by applying our algorithm to the liquid biopsy data of (Best et al., 2015). These data contain $p=58\,037$ genes sequenced in n=288 human subjects, 60 of whom were healthy controls, the others having been previously diagnosed with cancer. Averaging over 20 random 80/20 train/test splits of this dataset, we found that CoDaCoRe achieved the same predictive accuracy as competing methods (within error), but obtained a much sparser model. Remarkably, CoDaCoRe identified log-ratios involving just 3 genes, that were equally predictive to both black-box classifiers and linear models with over 20 active variables. This case study again illustrates the potential of CoDaCoRe to derive novel biological insights, and also to develop learning algorithms for cancer diagnosis, a domain in which model interpretability—including sparsity—is of paramount importance (Wan et al., 2017).

4.4 Simulation study

In addition to the above previous experiments, we provide a simulation study in Section G of Supplementary Material. For simulated HTS datasets of dimensionality ranging from 100 to 10 000 input variables, we find that CoDaCoRe is able to recover the true biomarkers used in the data-generating process, and does so with similar or higher accuracy (and orders of magnitude faster) than its competitors.

5 Conclusion

Our results corroborate the summary in Table 1: CoDaCoRe is the first sparse and interpretable CoDa model that can scale to high-dimensional HTS data. It does so convincingly, with linear scaling that results in runtimes similar to linear models. Our method is also

competitive in terms of predictive accuracy, performing comparably to powerful black-box classifiers, but with interpretability. Our findings suggest that CoDaCoRe could play a significant role in the future analysis of high-throughput sequencing data, with broad implications in microbiology, statistical genetics and the field of CoDa.

Funding

The authors thank the Simons Foundation 542963, Sloan Foundation, McKnight Endowment Fund, NSF DBI-1707398, and the Gatsby Charitable Foundation for support.

Conflict of Interest: none declared.

Data availability

The 25 benchmark datasets of Section 4.1 can be found at https://zenodo.org/record/3893986, and the simulated datasets together with further instructions can be found at our repo https://github.com/cunningham-lab/codacore. As for the liquid biopsy data of Section 4.3, we refer to the original publication (Best *et al.*, 2015).

References

Aitchison, J. (1982) The statistical analysis of compositional data. J. R. Stat. Soc. Ser. B (Methodological), 44, 139–160.

Alix-Panabières, C. and Pantel, K. (2016) Clinical applications of circulating tumor cells and circulating tumor DNA as liquid biopsy. *Cancer Discov.*, 6, 479–491.

Bates, S. and Tibshirani, R. (2019) Log-ratio lasso: scalable, sparse estimation for log-ratio models. *Biometrics*, 75, 613–624.

Best, M.G. et al. (2015) RNA-seq of tumor-educated platelets enables blood-based pan-cancer, multiclass, and molecular pathway cancer diagnostics. Cancer Cell, 28, 666–676.

Calle, M.L. (2019) Statistical analysis of metagenomics data. *Genomics Inf.*, 17, e6.

Cammarota, G. et al. (2020) Gut microbiome, big data and machine learning to promote precision medicine for cancer. Nat. Rev. Gastroenterol. Hepatol., 17, 635–648.

Crovesy, L. et al. (2020) Profile of the gut microbiota of adults with obesity: a systematic review. Eur. J. Clin. Nutr., 74, 1251–1262.

Dillies, M.-A. et al.; French StatOmique Consortium. (2013) A comprehensive evaluation of normalization methods for illumina high-throughput RNA sequencing data analysis. Brief. Bioinf., 14, 671–683.

Egozcue, J.J. and Pawlowsky-Glahn, V. (2005) Groups of parts and their balances in compositional data analysis. *Math. Geol.*, 37, 795–828.

Egozcue, J.J. (2003) Isometric logratio transformations for compositional data analysis. Math. Geol., 35, 279–300.

Fernandes, A.D. et al. (2013) Anova-like differential expression (ALDEX) analysis for mixed population RNA-seq. PLoS One, 8, e67019.

Fernandes, A.D. et al. (2014) Unifying the analysis of high-throughput sequencing datasets: characterizing RNA-seq, 16s RRNA gene sequencing and selective growth experiments by compositional data analysis. *Microbiome*, 2, 15

Filzmoser, P. and Walczak, B. (2014) What can go wrong at the data normalization step for identification of biomarkers? *J. Chromatography A*, 1362, 194–205.

Filzmoser, P. et al. (2009) Univariate statistical analysis of environmental (compositional) data: problems and possibilities. Sci. Total Environ., 407, 6100–6108.

Friedman, J. et al. (2001) The Elements of Statistical Learning, Vol. 1. Springer Series in Statistics, New York.

Gloor,G.B. and Reid,G. (2016) Compositional analysis: a valid approach to analyze microbiome high-throughput sequencing data. *Can. J. Microbiol.*, 62, 692–703.

Gloor, G.B. et al. (2016) It's all relative: analyzing microbiome data as compositions. Ann. Epidemiol., 26, 322–329.

Gloor, G.B. et al. (2017) Microbiome datasets are compositional: and this is not optional. Front. Microbiol., 8, 2224.

Goodman,B. and Flaxman,S. (2017) European union regulations on algorithmic decision-making and a "right to explanation". AI Mag., 38, 50–57.

- Greenacre, M. (2019a) Comments on: compositional data: the sample space and its structure. TEST, 28, 644–652.
- Greenacre, M. (2019b) Variable selection in compositional data analysis using pairwise logratios. *Math. Geosci.*, 51, 649–682.
- Greenacre, M. (2020) Amalgamations are valid in compositional data analysis, can be used in agglomerative clustering, and their logratios have an inverse transformation. *Appl. Comput. Geosci.*, 5, 100017.
- Greenacre, M. et al. (2020) A comparison of isometric and amalgamation logratio balances in compositional data analysis. Computers & Geosciences, 104, 104621
- He, H. and Ma, Y. (2013) Imbalanced Learning: Foundations, Algorithms, and Applications, Wiley-IEEE Press, New York.
- Jang, E. et al. (2017) Categorical reparameterization with gumbel-softmax. In: 5th International Conference on Learning Representations, (ICLR) 2017, Toulon. France. April 24-26. 2017. Conference Track Proceedings.
- Li,H. (2015) Microbiome, metagenomics, and high-dimensional compositional data analysis. Annu. Rev. Stat. Appl., 2, 73–94.
- Linderman, S. et al. (2018) Reparameterizing the Birkhoff polytope for variational permutation inference. In: International Conference on Artificial Intelligence and Statistics, (AISTATS) 2018, 9-11 April 2018, Playa Blanca, Lanzarote, Canary Islands, Spain.
- Lovell, D. et al. (2015) Proportionality: a valid alternative to correlation for relative data. PLoS Comput. Biol., 11, e1004075.
- Lu, J. et al. (2019) Generalized linear models with linear constraints for microbiome compositional data. Biometrics, 75, 235–244.
- Maddison, C. J. et al. (2017) The concrete distribution: a continuous relaxation of discrete random variables. In: 5th International Conference on Learning Representations, (ICLR) 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings.
- Magne, F. et al. (2020) The firmicutes/bacteroidetes ratio: a relevant marker of gut dysbiosis in obese patients? Nutrients, 12, 1474.
- Mena,G. et al. (2018) Learning latent permutations with Gumbel-Sinkhorn networks. In: 6th International Conference on Learning Representations, (ICLR) 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings.
- Morton, J.T. et al. (2017) Balance trees reveal microbial niche differentiation. MSystems, 2, e00162-16.
- Morton, J.T. et al. (2019) Establishing microbial composition measurement standards with reference frames. Nat. Commun., 10, 2719.
- Pawlowsky-Glahn, V. and Buccianti, A. (2011) Compositional Data Analysis: Theory and Applications. Wiley-Blackwell, Chichester, UK.
- Pawlowsky-Glahn, V. and Egozcue, J.J. (2006) Compositional data and their analysis: an introduction. Geol. Soc. Lond. Special Public., 264, 1–10.

- Pearson,K. (1896) VII. Mathematical contributions to the theory of evolution. III. Regression, heredity, and panmixia. *Philos. Trans. R. Soc. Lond. Ser. A*, 187, 253–318.
- Potapczynski, A. et al. (2020) Invertible gaussian reparameterization: revisiting the gumbel-softmax. Advances in Neural Information Processing Systems, 33.
- Prifti, E. et al. (2020) Interpretable and accurate prediction models for metagenomics data. GigaScience, 9, giaa010.
- Quinn, T. et al. (2020) Deepcoda: personalized interpretability for compositional health data. In: Proceedings of the 37th International Conference on Machine Learning, (ICML) 2020, 13-18 July 2020, Virtual Event.
- Quinn, T.P. and Erb, I. (2019) Using balances to engineer features for the classification of health biomarkers: a new approach to balance selection. *bioRxiv*, 600122.
- Quinn, T.P. and Erb,I. (2020) Amalgams: data-driven amalgamation for the dimensionality reduction of compositional data. NAR Genomics Bioinf., 2, 1qaa076.
- Quinn, T.P. et al. (2017) propr: an r-package for identifying proportionally abundant features using compositional data analysis. Sci. Rep., 7, 16252–16259.
- Quinn, T.P. et al. (2018) Understanding sequencing data as compositions: an outlook and review. Bioinformatics, 34, 2870–2878.
- Quinn, T.P. et al. (2019) A field guide for the compositional analysis of any-omics data. GigaScience, 8, giz107.
- Quinn, T.P. et al. (2021) A critique of differential abundance analysis, and advocacy for an alternative. arXiv, preprint arXiv:2104.07266.
- Rahat-Rozenbloom, S. et al. (2014) Evidence for greater production of colonic short-chain fatty acids in overweight than lean humans. Int. J. Obesity, 38, 1525–1531
- Rivera-Pinto, J. et al. (2018) Balances: a new perspective for microbiome analysis. MSystems. 3. e00053-18.
- Sheng, M. et al. (2018) Identification of tumor-educated platelet biomarkers of non-small-cell lung cancer. Onco Targets Ther., 11, 8143–8151.
- Silverman, J.D. et al. (2017) A phylogenetic transform enhances analysis of compositional microbiota data. Elife, 6, e21887.
- Susin, A. et al. (2020) Variable selection in microbiome compositional data analysis. NAR Genomics and Bioinformatics, 2, 1qaa029.
- Vangay,P. et al. (2019) Microbiome Learning Repo (ML Repo): a public repository of microbiome regression and classification tasks. GigaScience. 8.
- Wan, J.C. et al. (2017) Liquid biopsies come of age: towards implementation of circulating tumour DNA. Nat. Rev. Cancer, 17, 223–238.
- Washburne, A.D. et al. (2017) Phylogenetic factorization of compositional data yields lineage-level associations in microbiome datasets. Peer J. 5, e2969.
- Zhang, Y.-H. *et al.* (2017) Identifying and analyzing different cancer subtypes using RNA-seq data of blood platelets. *Oncotarget*, 8, 87494–87511.