# Multiple Augmented Reduced Rank Regression for Pan-Cancer Analysis

Jiuzhou Wang, Eric F. Lock

Division of Biostatistics, University of Minnesota, Minneapolis, MN, USA

September 1, 2023

## Abstract

Statistical approaches that successfully combine multiple datasets are more powerful, efficient, and scientifically informative than separate analyses. To address variation architectures correctly and comprehensively for high-dimensional data across multiple sample sets (i.e., cohorts), we propose multiple augmented reduced rank regression (maRRR), a flexible matrix regression and factorization method to concurrently learn both covariate-driven and auxiliary structured variation. We consider a structured nuclear norm objective that is motivated by random matrix theory, in which the regression or factorization terms may be shared or specific to any number of cohorts. Our framework subsumes several existing methods, such as reduced rank regression and unsupervised multi-matrix factorization approaches, and includes a promising novel approach to regression and factorization of a single dataset (aRRR) as a special case. Simulations demonstrate substantial gains in power from combining multiple datasets, and from parsimoniously accounting for all structured variation. We apply maRRR to gene expression data from multiple cancer types (i.e., pan-cancer) from TCGA, with somatic mutations as covariates. The method performs well with respect to prediction and imputation of held-out data, and provides new insights into mutation-driven and auxiliary variation that is shared or specific to certain cancer types.

**Keywords:** cancer genomics, data integration, low rank matrix factorization, missing data imputation, nuclear norm, reduced rank regression

# 1   Introduction

The proliferation of omics data in biomedicine and genomics has allowed for increasingly comprehensive investigations that span multiple sample sets and multiple molecular facets. Statistical approaches that successfully combine multiple datasets within a single analytical framework are more powerful, efficient, and scientifically informative than separate analyses. This has spurred advances in methodology for high-dimensional data integration, however, there remain unmet needs especially for multi-cohort data in which the same features are measured for different sample groups. Our motivating example is gene expression and somatic mutation data from the Cancer Genome Atlas (TCGA) Pan-Cancer Project (Hoadley et al., 2018; Hutter and Zenklusen, 2018), for 6581 tumor samples from 30 cohorts corresponding to different cancer types. Given the importance of gene expression in the behavior of cancer, and the related etiology of distinct cancer types through somatic mutations, we are interested in distinguishing variation due to somatic mutations from auxiliary structured variation in cancer gene expression and whether these effects are shared across cancer types.

Several unsupervised multi-matrix factorization methods provide low-rank representations of underlying structure. The singular value decomposition (SVD), principle component analysis (PCA) and other well-known approaches allow a rank $r$ approximation of a single matrix $\mathbf{X}_{p\times n} \approx \mathbf{U}_{p\times r}\mathbf{V}_{n\times r}^T, r < \min(p, n)$. Loadings $\mathbf{U}$ and scores $\mathbf{V}$ explain variation in the rows or columns, respectively. The joint and individual variation explained (JIVE) method extends PCA to multiple datasets with shared columns $\{\mathbf{X}_1, \ldots, \mathbf{X}_J\}$ via $\mathbf{X}_i \approx \mathbf{U}_i\mathbf{V}^T + \mathbf{W}_i\mathbf{V}_i^T$. Here the joint scores $\mathbf{V}$ capture shared structure among the datasets, and the individual scores $\mathbf{V}_i$ capture structure specific to dataset $i$. Numerous related approaches, such as AJIVE (Feng et al., 2018) and SLIDE (Gaynanova and Li, 2019) have been proposed to factorize multiple data from other perspectives. Moreover, BIDIFAC+ (Lock et al., 2022) enables a more flexible way to identify multiple shared and specific modules of variation, which may be partially shared over row subset or column subsets. However, these unsupervised methods suffer from neglecting covariate information. Other supervised techniques (Wang and Safo, 2021; Zhang and Gaynanova, 2022) identify structures across multiple datasets relevant to predicting an outcome, but they do not capture both covariate-driven and auxiliary structures.

To impose low-rank covariate effects, different types of penalties have been introduced in the the multivariate least square regression framework. Reduced rank regression (RRR) (Izenman, 1975) is a popular approach to predict $\mathbf{X} : p \times n$ from $\mathbf{Y} : q \times n$ via least squares in which the coefficients have low-rank, $\mathbf{X} \approx \mathbf{BY}$ with rank$(\mathbf{B}) < \min(p, q)$. Rank penalized (RSC) (Bunea et al., 2011) and nuclear-norm penalized (NNP) least square criteria (Yuan et al., 2007) are widely used alternatives with penalties that enforce low-rank coefficients. Combining RRR with adaptive NNP (Chen et al., 2013) shows a better performance than RSC. Integrative RRR (Li et al., 2019) extends the estimation to multiple covariate sets all at once. Nonetheless, those regression methods have two limitations: (1) they do not allow for potentially unique covariate-driven signals across multiple sample cohorts and (2) they do not account for additional low-rank structure unrelated to the covariates.

Missing values occur in genomics and other fields due to cost limitations or other technical issues. The data may have three types of missingness: entry-wise, column-wise or row-wise. To impute missing values matrix factorization based approaches, such as SVDImpute (Troyanskaya et al., 2001) and SoftImpute (Mazumder et al., 2010) are popular since they are effective and straightforward, and many of the the aforementioned methods can be modified for imputation. However, they will suffer from the same limitations described above.

Unifying reduced rank regression and unsupervised low-rank factorization using the nuclear norm penalty, we develop the multiple augmented reduced rank regression (maRRR) method for multi-cohort data that enables a very flexible approach for the simultaneous iden-

tification of covariate-driven effects and auxiliary structured variation. These covariate effects and augmented structures may be shared across any cohorts via a general objective function. This novel low-rank regression and factorization method can be used to impute various types of missing data, accurately capture the relationship between covariates and high-dimensional outcomes, and explore covariate-related and covariate-unrelated patterns of variation that are shared across or specific to different cohorts.

## 2 Proposed Model

Let $\mathbf{X}_j : p \times n_j$ denote data matrices with accompanying covariates $\mathbf{Y}_j : q \times n_j$ for $j$ sample cohorts $j = 1, ..., J$. Concatenations across all cohorts are denoted by $\cdot$, e.g., $\mathbf{X}_\cdot = [\mathbf{X}_1, \mathbf{X}_2, ..., \mathbf{X}_J]$ and $\mathbf{Y}_\cdot = [\mathbf{Y}_1, \mathbf{Y}_2, ..., \mathbf{Y}_J]$. Both $\mathbf{X}_\cdot$ and $\mathbf{Y}_\cdot$ are the only observed data in the model. For our application, we consider gene expression data $\mathbf{X}_\cdot$ and somatic mutations $\mathbf{Y}_\cdot$ for several patients across $J = 30$ cancer types. We are interested in decomposing $\mathbf{X}_\cdot$ into 'modules' of low-rank covariate-driven or auxiliary structures, where each module is shared on a different subset of the cohorts. We estimate low-rank coefficient matrices $\mathbf{B}_k : p \times q$ for $k = 1, ..., K$ modules of covariate-driven variation and we concurrently estimate low-rank auxiliary variation structures $\mathbf{S}_\cdot^{(l)} : p \times n$ for $l = 1, ..., L$ modules. Acknowledging the errors $\mathbf{E}_j : p \times n_j, j = 1, ..., J$ for each cohort, the full model is

$$\mathbf{X}_\cdot = \sum_{k=1}^{K} \mathbf{B}_k \mathbf{Y}_\cdot^{(k)} + \sum_{l=1}^{L} \mathbf{S}_\cdot^{(l)} + \mathbf{E}_\cdot \tag{1}$$

where $\mathbf{Y}_\cdot^{(k)} = [\mathbf{Y}_1^{(k)}, \mathbf{Y}_2^{(k)}, ..., \mathbf{Y}_J^{(k)}]$, $\mathbf{S}_\cdot^{(l)} = [\mathbf{S}_1^{(l)}, \mathbf{S}_2^{(l)}, ..., \mathbf{S}_J^{(l)}]$, $\mathbf{E}_\cdot = [\mathbf{E}_1, \mathbf{E}_2, ..., \mathbf{E}_J]$.

The presence of each $\mathbf{Y}_j^{(k)}$ or $\mathbf{S}_j^{(l)}$ across the cohorts are determined by binary indicator matrices $\mathbf{C}_Y : J \times K$ and $\mathbf{C}_S : J \times L$ respectively:

$$\mathbf{Y}_j^{(k)} = \begin{cases} \mathbf{0}_{q \times n_j} & \text{if } \mathbf{C}_Y[j, k] = 0 \\ \mathbf{Y}_j & \text{if } \mathbf{C}_Y[j, k] = 1, \end{cases} \quad \mathbf{S}_j^{(l)} = \begin{cases} \mathbf{0}_{p \times n_j} & \text{if } \mathbf{C}_S[j, l] = 0 \\ \mathbf{U}_S^{(l)} \mathbf{V}_{Sj}^{(l)T} & \text{if } \mathbf{C}_S[j, l] = 1. \end{cases}$$

$\mathbf{U}_S^{(l)}$ represents shared loadings and $\mathbf{V}_{Sj}^{(l)}$ sample scores for cohort $j$ in module $l$. Both indicator matrices may be determined either by pre-existing knowledge or via a data-driven algorithm, which we will detail in Section 7.2 and Appendix A.2. They are fixed in the model estimation process. There should be no identical columns within $\mathbf{C}_Y$, so that each $\mathbf{B}_k \mathbf{Y}_\cdot^{(k)}$ is present on a distinct subset of the cohorts. Similarly, no duplicate columns within $\mathbf{C}_S$. We refer to each $\mathbf{B}_k \mathbf{Y}_\cdot^{(k)}$ and $\mathbf{S}_\cdot^{(l)}$ as a module. Each $\mathbf{S}_\cdot^{(l)}$ gives a low-rank module that explains covariate-unrelated structured variability within the cohorts (e.g., cancer types) identified by $\mathbf{C}_S[:, l]$. Each $\mathbf{B}_k \mathbf{Y}_\cdot^{(k)}$ gives another low-rank module for covariate-driven structure for the cancer type identified by $\mathbf{C}_Y[:, l]$. Each module is assumed to be low-rank, meaning it can be factorized as the product of a small number of row and column vectors, $\mathbf{B}_k = \mathbf{U}_B^{(k)} \mathbf{V}_B^{(k)T}$ and $\mathbf{S}_\cdot^{(l)} = \mathbf{U}_S^{(l)} \mathbf{V}_S^{(l)T}$. We provide a schematic of our model in Fig. 1 and a table of notation details in Appendix A.1.
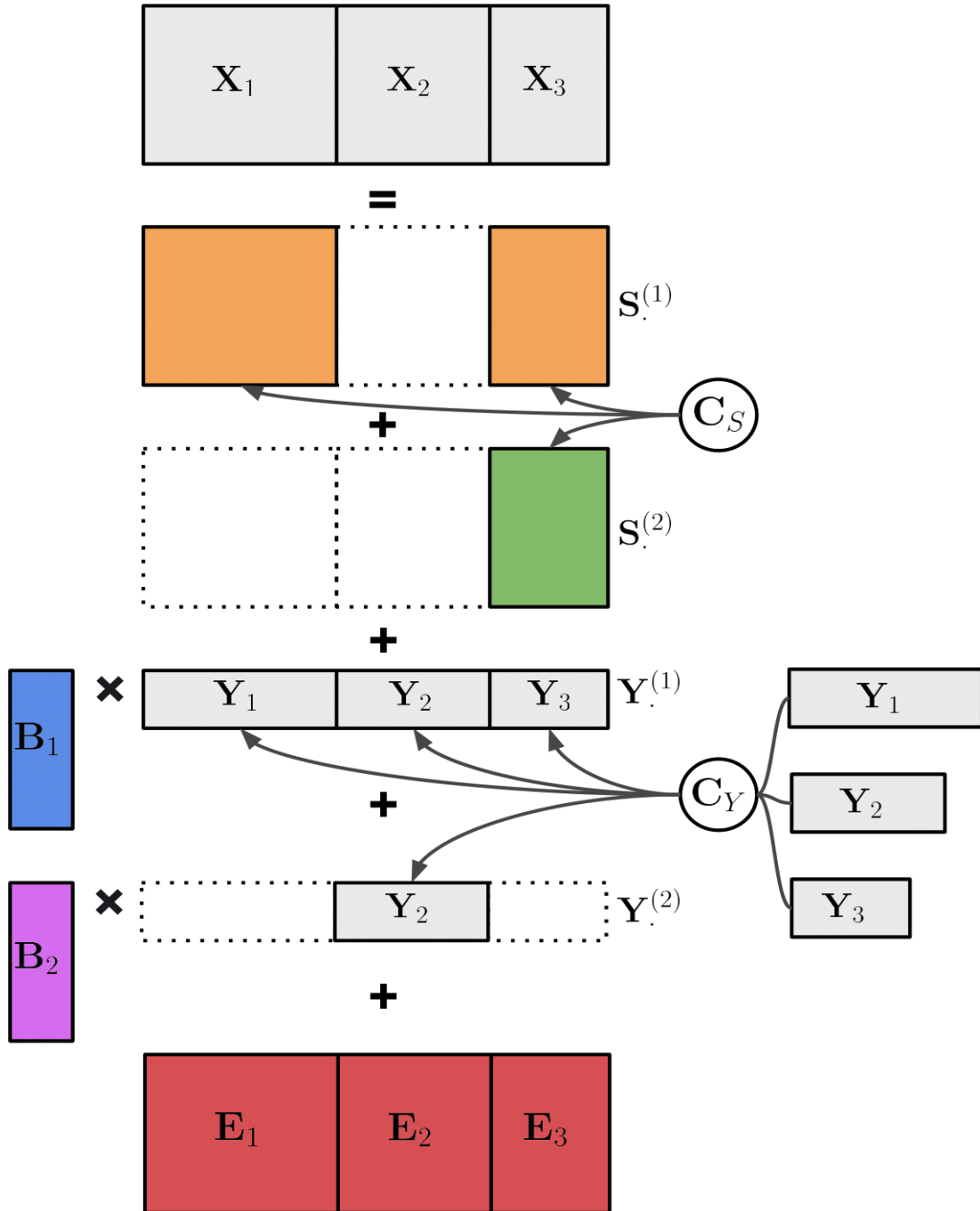
Figure 1: A schematic of our proposed model maRRR with 3 cohorts as an example. All matrices in grey are observed, i.e. outcomes $\mathbf{X}. = [\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3]$ and covariates $\mathbf{Y}. = [\mathbf{Y}_1, \mathbf{Y}_2, \mathbf{Y}_3]$. Two binary indicator matrices for auxiliary structures $\mathbf{C}_S = [[1, 0, 1]^T, [0, 0, 1]^T]$ and for covariate effects $\mathbf{C}_Y = [[1, 1, 1]^T, [0, 1, 0]^T]$ are pre-specified. Then, the structures of $\mathbf{S}^{(1)} = [\mathbf{U}_S^{(1)}\mathbf{V}_{S1}^{(1)T}, \mathbf{0}, \mathbf{U}_S^{(1)}\mathbf{V}_{S3}^{(1)T}]$, $\mathbf{S}^{(2)} = [\mathbf{0}, \mathbf{0}, \mathbf{U}_S^{(2)}\mathbf{V}_{S3}^{(2)T}]$ and $\mathbf{Y}^{(1)} = [\mathbf{Y}_1, \mathbf{Y}_2, \mathbf{Y}_3], \mathbf{Y}^{(2)} = [\mathbf{0}, \mathbf{Y}_2, \mathbf{0}]$ are determined. All matrices in color are to estimate, i.e., auxiliary structures $\mathbf{S}^{(1)}, \mathbf{S}^{(2)}$, covariate effect coefficients $\mathbf{B}_1, \mathbf{B}_2$ and random errors $\mathbf{E}. = [\mathbf{E}_1, \mathbf{E}_2, \mathbf{E}_3]$.

4

# 3 Objective Function

To estimate model (1) and impose low-rank structure, we minimize the following least squares criterion with a structured nuclear norm penalty:

$$\min_{\{\mathbf{B}_k\}_{k=1}^K, \{\mathbf{S}_.^{(l)}\}_{l=1}^L} \left\{ \frac{1}{2} ||\mathbf{X}_. - \sum_{k=1}^K \mathbf{B}_k \mathbf{Y}_.^{(k)} - \sum_{l=1}^L \mathbf{S}_.^{(l)}||_F^2 + \sum_{k=1}^K \lambda_B^{(k)} ||\mathbf{B}_k||_* + \sum_{l=1}^L \lambda_S^{(l)} ||\mathbf{S}_.^{(l)}||_* \right\} \quad (2)$$

Here $|| \cdot ||_*$ denotes the nuclear norm, i.e., the sum of the singular values of the matrix, a convex penalty which encourages a low-rank solution. There are three special cases of the general objective functions worth noting. The first two are novel and the third is previously described, listed as follows:

1. Augmented reduced rank regression (aRRR), our proposed approach minimizing (2) for $K = L = J = 1$. The nuclear-norm penalized reduced rank regression model for a single matrix is "augmented" to account for auxiliary structured variation $\mathbf{S}$ simultaneously.

2. Multi-cohort reduced rank regression (mRRR), our proposed approach minimizing (2) for $L = 0$, i.e. no auxiliary terms $\mathbf{S}$. The reduced rank regression is extended to recover multiple (shared or individual) covariate effects at once.

3. Optimizing this objective with no covariate-driven structure ($K = 0$) corresponds to the BIDIFAC+ method (Lock et al., 2022) with horizontal structures only.

Mazumder et al. (2010) and others have noted the equivalence of the nuclear norm penalty and an additive $L_2$ penalty on the terms in the low-rank factorization, and this leads to an alternative form of our objective (Eq. (2)),

$$\min_{\{\mathbf{U}_B^{(k)}, \mathbf{V}_B^{(k)}\}_{k=1}^K, \{\mathbf{U}_S^{(l)}, \mathbf{V}_S^{(l)}\}_{l=1}^L} \frac{1}{2} \left\{ ||\mathbf{X}_. - \sum_{k=1}^K \mathbf{U}_B^{(k)} \mathbf{V}_B^{(k)T} \mathbf{Y}_.^{(k)} - \sum_{l=1}^L \mathbf{U}_S^{(l)} \mathbf{V}_S^{(l)T}||_F^2 + \right.$$
$$\left. \sum_{k=1}^K \lambda_B^{(k)} (||\mathbf{U}_B^{(k)}||_F^2 + ||\mathbf{V}_B^{(k)}||_F^2) + \sum_{l=1}^L \lambda_S^{(l)} (||\mathbf{U}_S^{(l)}||_F^2 + ||\mathbf{V}_S^{(l)}||_F^2) \right\} \quad (3)$$

where we only need to set a general upper bound for the estimated rank of each $\mathbf{B}_k$ and $\mathbf{S}_.^{(l)}$, i.e. $r_{B,upper}$ and $r_{S,upper}$. These upper bounds serve as the number of columns for each $\mathbf{U}_B^{(k)}$, $\mathbf{V}_B^{(k)}$, $\mathbf{U}_S^{(l)}$, and $\mathbf{V}_S^{(l)}$. The actual ranks of the solution may be smaller due to the rank sparsity encouraged by the nuclear norm penalty, and if the upper bounds are large enough the solution will correspond to that in equation (2). We state this formally in Theorem 1; the proof of this result, and all other novel results in this manuscript, are given in Appendix B.

**Theorem 1.** *If both (2) and (3) have the same penalty terms $\lambda_B^{(k)} > 0, k = 1, ..., K$ and $\lambda_S^{(l)} > 0, l = 1, ..., L$, the solutions to the objective functions coincide.*

In what follows in Section 4 we describe a random matrix theory approach to automatically select the nuclear norm penalty weights $\lambda$ for the different modules.

# 4 Theoretical results

We describe conditions on the penalty to avoid degenerate cases in which certain modules are guaranteed to be zero in the solution (regardless of the data $\mathbf{X}_.$ and $\mathbf{Y}_.$) in Proposition 1.

**Proposition 1.** *The following conditions on penalty parameters are needed to allow for non-zero estimation of each $\{\mathbf{B}_k\}_{k=1}^{K}, \{\mathbf{S}^{(l)}\}_{l=1}^{L}$:*

1. *Let $\mathcal{I}_k \subset \{1, ..., k-1, k+1, ..., K\}$ be any subset of $\mathbf{Y}$ modules for which the non-zero blocks of $\{\mathbf{Y}_{\cdot}^{(i)}\}_{i\in\mathcal{I}_k}$ cover exactly those of $\mathbf{Y}^{(k)}$, i.e. $\sum_{i\in\mathcal{I}_k} \mathbf{C}_Y[\cdot, i] = c_y \cdot \mathbf{C}_y[\cdot, k]$ for some positive integer $c_y$. Then, $\lambda_B^{(k)} < \frac{1}{c_y}\sum_{i\in\mathcal{I}_k}\lambda_B^{(i)}$.*

2. *Let $\mathcal{I}_k \subset \{1, 2, ..., L\}$ be any subset of $\mathbf{S}$ modules for which the non-zero blocks of $\{\mathbf{S}_{\cdot}^{(i)}\}_{i\in\mathcal{I}_k}$ cover exactly those of $\mathbf{Y}_k$, i.e. $\sum_{i\in\mathcal{I}_k} \mathbf{C}_S[\cdot, i] = c_{sy} \cdot \mathbf{C}_Y[\cdot, k]$ for some positive integer $c_{sy}$. Then, $\lambda_B^{(k)} < \frac{1}{c_{sy}}\sum_{i\in\mathcal{I}_k}\lambda_S^{(i)}||\mathbf{Y}_{\cdot}^{(k)}||_*$.*

3. *For $l \neq l'$, if a module $\mathbf{S}_{\cdot}^{(l')}$ is contained in another module $\mathbf{S}_{\cdot}^{(l)}$, i.e. $\mathbf{C}_S[j, l] \geq \mathbf{C}_S[j, l'], \forall j$, then $\lambda_S^{(l')} < \lambda_S^{(l)}$.*

4. *Let $\mathcal{I}_l \subset \{1, ..., l-1, l+1, ..., L\}$ be any subset of $\mathbf{S}$ modules that the non-zero blocks of $\{\mathbf{S}_{\cdot}^{(i)}\}_{i\in\mathcal{I}_l}$ cover exactly those of $\mathbf{S}_{\cdot}^{(l)}$, i.e. $\sum_{i\in\mathcal{I}_l} \mathbf{C}_S[\cdot, i] = c_s \cdot \mathbf{C}_S[\cdot, l]$ for some positive integer $c_s$. Then, $\lambda_S^{(l)} < \frac{1}{c_s}\sum_{i\in\mathcal{I}_l}\lambda_S^{(i)}$.*

To motivate a random matrix theory approach to select the tuning parameters, we present two results establishing the connection between the nuclear norm penalty and singular value thresholding. Lemma 1 is a well-known result for the unsupervised case (Cai et al., 2010), and in Proposition 2 we extend it to the regression context.

**Lemma 1.** *Let $\mathbf{U}\mathbf{D}\mathbf{V}^T$ be the SVD of a matrix $\mathbf{X}$. The solution to $\min_{\mathbf{S}}\{\frac{1}{2}||\mathbf{X}-\mathbf{S}||_F^2 + \lambda||\mathbf{S}||_*\}$ is $\mathbf{S} = \mathbf{U}\widetilde{\mathbf{D}}\mathbf{V}^T$, where $\widetilde{\mathbf{D}}$ is diagonal with entries $\widetilde{\mathbf{D}}[i, i] = \max(\mathbf{D}[i, i] - \lambda, 0)$.*

**Proposition 2.** *Let $\mathbf{Y}$ be a semi-orthogonal matrix such that $\mathbf{Y}\mathbf{Y}^T = \mathbf{I}$ and $\mathbf{U}\mathbf{D}\mathbf{V}^T$ be the SVD of a matrix $\mathbf{X}\mathbf{Y}^T$. The solution to both of the following objectives:*

$$\min_{\mathbf{B}}\{\frac{1}{2}||\mathbf{X} - \mathbf{B}\mathbf{Y}||_F^2 + \lambda||\mathbf{B}||_*\} \quad and \quad \min_{\mathbf{B}}\{\frac{1}{2}||\mathbf{X} - \mathbf{B}\mathbf{Y}||_F^2 + \lambda||\mathbf{B}\mathbf{Y}||_*\},$$

*is $\mathbf{B} = \mathbf{U}\widetilde{\mathbf{D}}\mathbf{V}^T$, where $\widetilde{\mathbf{D}}$ is diagonal with entries $\widetilde{\mathbf{D}}[i, i] = \max(\mathbf{D}[i, i] - \lambda, 0)$.*

While the relative merits of penalizing $||\mathbf{B}||_*$ or $||\mathbf{B}\mathbf{Y}||_*$ has been debated (Yuan et al., 2007; Chen et al., 2013), Proposition 2 shows they are identical if $\mathbf{Y}$ is semi-orthogonal. In practice, we orthogonalize the columns of $\mathbf{Y}$ prior to estimation. However, this requires that the number of features in $\mathbf{Y}$ is less than the sample size (e.g., $q < n$); if $q \geq n$ then $\mathbf{Y}$ will be semi-orthogonal in the opposite direction $\mathbf{Y}^T\mathbf{Y} = \mathbf{I}$, causing the solution to degenerate to the unsupervised case, which we establish in Proposition 5 in Appendix B.

The following propositions describe the distribution of the singular values of a random matrix under general assumptions, which can then be used to motivate tuning parameters.

**Proposition 3.** *Let $\lambda_{max}$ be the largest singular value of a matrix $\mathbf{E} : m \times n$ of independent Guassian entries with mean 0 and variance $\sigma^2$. We have $E(\lambda_{max}) \leq \sigma(\sqrt{m} + \sqrt{n})$.*

**Proposition 4.** *Let $\mathbf{Y}_{q\times n}$ be semi-orthogonal such that $\mathbf{Y}\mathbf{Y}^T = \mathbf{I}$. For integers $m, q \geq 1$ defined in a way that $\frac{m}{q} \to c > 0$ as $q \to \infty$, Let $\mathbf{X}_{m\times n}, \mathbf{B}_{m\times q}, \mathbf{E}_{m\times n}$ be three matrices such that $\mathbf{X} = \mathbf{B}\mathbf{Y}_{q\times n} + \frac{1}{\sqrt{q}}\mathbf{E}$, where entries of $\mathbf{E}$ are independent Guassian with mean 0 and variance $\sigma^2$. Assume $rank(\mathbf{B}) = r$. Denote the singular values of $\mathbf{B}$ and $\mathbf{X}\mathbf{Y}^T$ are $\sigma_1(\mathbf{B}) \geq ... \geq \sigma_r(\mathbf{B}) > 0$ and $\sigma_1(\mathbf{X}\mathbf{Y}^T) \geq ... \geq \sigma_r(\mathbf{X}\mathbf{Y}^T) > 0$ respectively. As $n \to \infty$,*

$$\sigma_j(\mathbf{X}\mathbf{Y}^T) \xrightarrow{P} \begin{cases} s(\sigma_j(\mathbf{B})) > 1 + \sqrt{c}, & if \ \sigma_j(\mathbf{B}) > \sqrt[4]{c} \\ 1 + \sqrt{c}, & if \ \sigma_j(\mathbf{B}) \leq \sqrt[4]{c} \end{cases}, \forall 1 \leq j \leq r,$$

where $s(\cdot)$ is a known function. In particular, when $\mathbf{Y}$ is an identity matrix $(q = n)$ and $\mathbf{X} = \mathbf{B} + \frac{1}{\sqrt{n}}\mathbf{E}$, it follows that $\sigma_j(\mathbf{X}) \xrightarrow{P} \begin{cases} s(\sigma_j(\mathbf{B})) > 1 + \sqrt{c}, & \text{if } \sigma_j(\mathbf{B}) > \sqrt[4]{c} \\ 1 + \sqrt{c}, & \text{if } \sigma_j(\mathbf{B}) \leq \sqrt[4]{c}. \end{cases}$

Proposition 3 comes directly from (Rudelson and Vershynin, 2010), and Proposition 4 is closely related to the result in (Shabalin and Nobel, 2013). Consider the reasonable penalty for $\mathbf{S}$ in Lemma 1, i.e. $\mathbf{X}_{m \times n} = \mathbf{S}_{m \times n} + \mathbf{E}_{m \times n}$. A set of reasonable tuning parameters will (1) detect the low-rank signals and (2) not capture components that are solely due to noise. Considering Propositions 3 and 4, setting $\lambda = \sigma(\sqrt{m} + \sqrt{n})$ is reasonable because it only keeps the signals (top $r$ components) whose singular values are expected to be greater than those of independent random noise. Consider the reasonable penalty for $\mathbf{B}$ in Proposition 2, i.e. $\mathbf{X}_{m \times n} = \mathbf{B}_{m \times q}\mathbf{Y}_{q \times n} + \mathbf{E}_{m \times n}$. Following a similar argument, we set $\lambda = \sigma(\sqrt{m} + \sqrt{q})$.

In practice, after normalizing raw data as described in Appendix C, the noise variance for $\mathbf{X}$ is 1 $(\sigma = 1)$ and each $\mathbf{Y}^{(k)}$ are semi-orthogonal. Thus, in order to distinguish true signals $\{\mathbf{B}_k\}_{k=1}^{K}, \{\mathbf{S}^{(l)}\}_{l=1}^{L}$ from Gaussian noise in the objective (2), we fix $\lambda_B^{(k)} = \sqrt{p} + \sqrt{q}$ for any module $\mathbf{B}_k, k = 1, ..., K$ and $\lambda_S^{(l)} = \sqrt{p} + \sqrt{\sum_{j=1}^{J} n_j \mathbf{C}_S[j, l]}$ for any module $\mathbf{S}^{(l)}, l = 1, ..., L$. This directly extends our choices for a single matrix, as estimating any given module $\mathbf{B}^{(k)}$ or $\mathbf{S}^{(l)}$ with the others fixed reduces to the setting of the previous propositions.

## 5  Estimation

In practice we scale $\mathbf{X}$ (Gavish and Donoho, 2017) and orthogonalize $\mathbf{Y}$ prior to optimization. The details of this procedure are provided in Appendix C, and a simulation study illustrating its advantages is provided in Appendix D.3.

### 5.1  Optimization

We estimate all regression coefficients $\mathbf{B}$ and auxiliary variation sources $\mathbf{S}$ simultaneously, via alternating optimization approaches for either formulation (2) or (3) of our objective. For objective (3), the introduction of $\mathbf{U}$ and $\mathbf{V}$ can make the optimization algorithm more efficient because the objective function has a closed-form gradient. Given all other estimates, we update every single $\mathbf{U}_B^{(k)}, \mathbf{V}_B^{(k)}, \mathbf{U}_S^{(l)}, \mathbf{V}_S^{(l)}$ by setting its corresponding gradient to be zero. The details are provided in Algorithm 1.

The symbol $\bigotimes$ means Kronecker product.

Note that Algorithm 1 does not require the columns of $\mathbf{Y}^{(k)}$ to be orthogonal. When $\mathbf{Y}^{(k)}$ is semi-orthogonal, in light of Lemma 1 and Proposition 2, we develop an alternative approach based on iterative soft-singular value thresholding estimators for (2) in Algorithm 2.

---
**Algorithm 1** Alternating Least Square with Matrix Decomposition
---
**Input:** Covariates $\mathbf{Y}$ and corresponding multivariate outcomes $\mathbf{X}$; penalizing terms $\lambda_B, \lambda_S$; binary indicator matrices $\mathbf{C}_Y, \mathbf{C}_S$

**Output:** $\mathbf{B}, \mathbf{S}$

1: **Initialization** Construct $\{\mathbf{Y}_{\cdot}^{(k)}\}_{k=1}^{K}$ based on $\mathbf{C}_Y$. Assign initialized numbers for each entry of $\{\mathbf{U}_B^{(k)}, \mathbf{V}_B^{(k)}\}_{k=1}^{K}, \{\mathbf{U}_S^{(l)}, \mathbf{V}_S^{(l)}\}_{l=1}^{L}$

2: **while** convergence criterion does not meet **do**

3:   **for** $k = 1, ..., K$ **do**

4:     Compute the residual matrix $\mathbf{X}_{\cdot}^{(k)} = \mathbf{X}_{\cdot} - \sum_{k'=1, k' \neq k}^{K} \mathbf{U}_B^{(k')} \mathbf{V}_B^{(k')T} \mathbf{Y}_{\cdot}^{(k')} - \sum_{l=1}^{L} \mathbf{U}_S^{(l)} \mathbf{V}_S^{(l)T}$

5:     Update $\mathbf{U}_B^{(k)} = \mathbf{X}_{\cdot}^{(k)} \mathbf{Y}_{\cdot}^{(k)T} \mathbf{V}_B^{(k)} (\mathbf{V}_B^{(k)T} \mathbf{Y}_{\cdot}^{(k)} \mathbf{Y}_{\cdot}^{(k)T} \mathbf{V}_B^{(k)} + \lambda_B^{(k)} \mathbf{I}_{r_B})^{-1}$

6:     Update $vec(\mathbf{V}_B^{(k)}) = [(\mathbf{U}_B^{(k)T} \mathbf{U}_B^{(k)}) \bigotimes (\mathbf{Y}_{\cdot}^{(k)} \mathbf{Y}_{\cdot}^{(k)T}) + \lambda_B^{(k)} \mathbf{I}_{q*r_B}]^{-1} vec[\mathbf{Y}_{\cdot}^{(k)} (\mathbf{X}_{\cdot}^{(k)T}) \mathbf{U}_B^{(k)}]$

7:     Transform $vec(\mathbf{V}_B^{(k)})$ to $\mathbf{V}_B^{(k)}$

8:   **end for**

9:   **for** $l = 1, .., L$ **do**

10:     Compute the residual matrix $\mathbf{X}_{\cdot}^{(l)} = \mathbf{X}_{\cdot} - \sum_{k=1}^{K} \mathbf{U}_B^{(k)} \mathbf{V}_B^{(k)T} \mathbf{Y}_{\cdot}^{(k)} - \sum_{l'=1, l' \neq l}^{L} \mathbf{U}_S^{(l')} \mathbf{V}_S^{(l')T}$

11:     Set $\mathbf{X}_j^{(l)} = 0$ where $\mathbf{C}_s[j, l] = 0$ for $j = 1, ..., J$

12:     Update $\mathbf{U}_S^{(l)} = \mathbf{X}_{\cdot}^{(l)} \mathbf{V}_S^{(l)} (\mathbf{V}_S^{(l)T} \mathbf{V}_S^{(l)} + \lambda_S^{(l)} \mathbf{I}_{r_S})^{-1}$

13:     Update $\mathbf{V}_S^{(l)} = \mathbf{X}_{\cdot}^{(l)T} \mathbf{U}_S^{(l)} (\mathbf{U}_S^{(l)T} \mathbf{U}_S^{(l)} + \lambda_S^{(l)} \mathbf{I}_{r_S})^{-1}$

14:   **end for**

15: **end while**

16: Set $\mathbf{B}_k = \mathbf{U}_B^{(k)} \mathbf{V}_B^{(k)T}$ for all $k = 1, .., K$, and $\mathbf{S}_{\cdot}^{(l)} = \mathbf{U}_S^{(l)} \mathbf{V}_S^{(l)T}$ for all $l = 1, .., L$

---

---
**Algorithm 2** Alternating Least Square with Soft-threshold Estimators
---
**Input:** Orthogonal covariates $\mathbf{Y}$ and corresponding multivariate outcomes $\mathbf{X}$; penalizing terms $\lambda_B, \lambda_S$; binary indicator matrices $\mathbf{C}_Y, \mathbf{C}_S$

**Output:** $\mathbf{B}, \mathbf{S}$

1: **Initialization** Construct $\{\mathbf{Y}_{\cdot}^{(k)}\}_{k=1}^{K}$ based on $\mathbf{C}_Y$. Assign initialized numbers for each entry of $\{\mathbf{B}_k\}_{k=1}^{K}, \{\mathbf{S}_{\cdot}^{(l)}\}_{l=1}^{L}$

2: **while** convergence criterion does not meet **do**

3:   **for** $k = 1, .., K$ **do**

4:     Compute the residual matrix $\mathbf{X}_{\cdot}^{(k)} = \mathbf{X}_{\cdot} - \sum_{k'=1, k' \neq k}^{K} \mathbf{B}_{k'} \mathbf{Y}_{\cdot}^{(k')} - \sum_{l=1}^{L} \mathbf{S}_{\cdot}^{(l)}$

5:     Compute the SVD of $\mathbf{X}_{\cdot}^{(k)} \mathbf{Y}_{\cdot}^{(k)T}$, i.e. $\mathbf{X}_{\cdot}^{(k)} \mathbf{Y}_{\cdot}^{(k)T} = \mathbf{L}_B^{(k)} \mathbf{D}_B^{(k)} \mathbf{R}_B^{(k)}$

6:     Update $\mathbf{B}_k = \mathbf{L}_B^{(k)} \widehat{\mathbf{D}}_B^{(k)} \mathbf{R}_B^{(k)}$ where $\widehat{\mathbf{D}}_B^{(k)}$ is a diagonal matrix with $\widehat{\mathbf{D}}_B^{(k)}[r, r] = max(\mathbf{D}_B^{(k)}[r, r] - \lambda_B^{(k)}, 0)$ for $r = 1, 2, ...$ on its diagonal entries and zero otherwise

7:   **end for**

8:   **for** $l = 1, .., L$ **do**

9:     Compute the residual matrix $\mathbf{X}_{\cdot}^{(l)} = \mathbf{X}_{\cdot} - \sum_{k=1}^{K} \mathbf{B}_{k'} \mathbf{Y}_{\cdot}^{(k')} - \sum_{l=1, l' \neq l}^{L} \mathbf{S}_{\cdot}^{(l')}$

10:     Set $\mathbf{X}_j^{(l)} = 0$ where $\mathbf{C}_s[j, l] = 0$ for $j = 1, ..., J$

11:     Compute the SVD of $\mathbf{X}_{\cdot}^{(l)}$, i.e. $\mathbf{X}_{\cdot}^{(l)} = \mathbf{L}_S^{(l)} \mathbf{D}_S^{(l)} \mathbf{R}_S^{(l)}$

12:     Update $\mathbf{S}_{\cdot}^{(l)} = \mathbf{L}_S^{(l)} \widehat{\mathbf{D}}_S^{(l)} \mathbf{R}_S^{(l)}$ where $\widehat{\mathbf{D}}_S^{(l)}$ is a diagonal matrix with $\widehat{\mathbf{D}}_S^{(l)}[r, r] = max(\mathbf{D}_S^{(l)}[r, r] - \lambda_S^{(l)}, 0)$ for $r = 1, 2, ...$ on its diagonal entries and zero otherwise

13:   **end for**

14: **end while**

---

For both algorithms, we use the same convergence criteria to decide whether to stop the optimization process: $\sum_{k=1}^{K} ||\widehat{\mathbf{B}}_k - \widetilde{\mathbf{B}}_k||_F^2 + \sum_{l=1}^{L} ||\widehat{\mathbf{S}}_{\cdot}^{(l)} - \widetilde{\mathbf{S}}_{\cdot}^{(l)}||_F^2 < \epsilon$, where $\widehat{\phantom{x}}$ denotes the estimation in the current epoch and $\widetilde{\phantom{x}}$ denotes the estimation in the previous epoch. It is also reasonable to use convergence of the loss function as the criteria.

In practice both algorithms have different strengths and weaknesses. Theoretically, Algorithm 2 can be used only when we orthogonalize the original $\mathbf{Y}$, because otherwise soft-thresholding to update $\mathbf{B}$ is not possible. In general, we find that the algorithms require similar computation time to achieve the same convergence criterion: Algorithm 1 tends to require less time if the true ranks (and accompanying maximum ranks specified, i.e. $r_{B,upper}$ and $r_{S,upper}$) are small, while Algorithm 2 is quicker and consumes less computational resources when the true rank and maximum ranks specific for Algorithm 1 are large.

## 5.2 Missing data imputation

One of the main uses of our proposed method is to impute various types of missing data. Based on the assumption that the abundance of existing entries provides sufficient information to uncover the global structures (both covariate and auxiliary effects) and therefore, to estimate the values of absent entries. Denote the set of indexes of all missing entries as $M$. Our iterative imputation process is as follows: (1) Initialize $\widetilde{\mathbf{X}}_{\cdot}$ by $\widetilde{\mathbf{X}}_{\cdot}[m,n] = \mathbf{X}_{\cdot}[m,n]$ if $[m,n] \notin M$, otherwise 0; (2) Estimate $\{\mathbf{B}_k\}_{k=1}^{K}$, $\{\mathbf{S}_{\cdot}^{(l)}\}_{l=1}^{L}$ by Algorithm 1 or 2 with current $\widetilde{\mathbf{X}}_{\cdot}$; (3) Update $\widetilde{\mathbf{X}}_{\cdot}$ by setting $\widetilde{\mathbf{X}}_{\cdot}[m,n] = (\sum_{k=1}^{K} \mathbf{B}_k \mathbf{Y}_{\cdot}^{(k)} + \sum_{l=1}^{L} \mathbf{S}_{\cdot}^{(l)})[m,n]$ for all $[m,n] \in M$; (4) Back to (2) unless convergence; the final $\widetilde{\mathbf{X}}_{\cdot}$ is the imputation result. This can be considered a modified EM-algorithm, and is similar to the approach used for softImpute (Mazumder et al., 2010) for nuclear-norm penalized imputation of a single matrix with no covariates.

# 6 Simulations

## 6.1 Recovery of true structure for special cases

Here, we present simulations as proof-of-concept for two novel scenarios within our approach: (i) simultaneous modeling of covariate effects and auxiliary low-rank variation and (ii) simultaneous modeling of shared or specific covariate effects across multiple cohorts.

For (i), we consider a single data matrix $\mathbf{X} : 100 \times 100$ and single set of covariates $\mathbf{Y} : 10 \times 100$ and generate data via $\mathbf{X} = \mathbf{BY} + \mathbf{S} + \mathbf{E}$, where $\mathbf{BY}$ is covariate-driven variation, $\mathbf{S}$ is auxiliary structured variation, and $\mathbf{E}$ is error. The coefficient array $\mathbf{B}$ has rank $R_y$ via $\mathbf{B} = a\mathbf{U}_B\mathbf{V}_B^T$ where $\mathbf{U}_B : 100 \times R_y$ and $\mathbf{V}_B : 10 \times R_y$, and $\mathbf{S}$ has rank 5 via $\mathbf{S} = b\mathbf{U}_S\mathbf{V}_S$ where $\mathbf{U}_S : 100 \times 5$ and $\mathbf{V}_S : 5 \times 100$. The entries of $\mathbf{E}$, $\mathbf{Y}$, $\mathbf{U}_B$, $\mathbf{V}_B$, $\mathbf{U}_S$ and $\mathbf{V}_S$ are all generated independently from a standard normal distribution. We consider $R_y = 1$ or $R_y = 5$, and consider three conditions with different signal strength for each term by adjusting $a$ and $b$: sd($\mathbf{BY}$) = 0.5 and sd($\mathbf{S}$) = 5 ($||\mathbf{BY}||/||\mathbf{S}|| = 0.1$, sd($\mathbf{BY}$) = sd($\mathbf{S}$) = 1 ($||\mathbf{BY}||/||\mathbf{S}|| = 1$), and sd($\mathbf{BY}$) = 5 and sd($\mathbf{S}$) = 0.5 ($||\mathbf{BY}||/||\mathbf{S}|| = 10$). For each set of conditions, we estimate $\mathbf{B}$ and $\mathbf{S}$ using four approaches: (1) Augmented reduced rank regression (aRRR), our proposed approach as described in Section 3, given 10 as the rank upper bound for $\mathbf{B}$ and $\mathbf{S}$. (2) Supervised singular value decomposition (SupSVD) (Li et al., 2016), a related model of the form $\mathbf{X} = \mathbf{YBV}^T + \mathbf{FV}^T + \mathbf{E}$ for one cohort, where $\mathbf{F}$ is the matrix of latent variables that correspond to auxiliary variation not related to the covariates, estimated using maximum likelihood and given the true rank of $\mathbf{BV}^T$ and $\mathbf{FV}^T$. (3) Two-stage least squares, in which the coefficients $\mathbf{B}$ is determined by ordinary least squares regression and $\mathbf{S}$ is determined by an SVD approximation with the true rank ($R = 5$) on the residuals $\mathbf{X} - \hat{\mathbf{B}}\mathbf{Y}$. (4) Two-stage nuclear norm (NN), in which $\mathbf{B}$ is determined by an NN-penalized reduced rank regression

and $\mathbf{S}$ by a NN-penalized matrix approximation to the residuals $\mathbf{X} - \hat{\mathbf{B}}\mathbf{Y}$. For each method, we compute the relative mean squared error (MSE) for $\mathbf{B}$ and $\mathbf{S}$, e.g., $||\mathbf{B} - \hat{\mathbf{B}}||_F^2/||\mathbf{B}||_F^2$. Average relative MSEs for each condition, over 100 replications, are shown in Table 1A. This demonstrates clear advantages of a nuclear norm penalty on $\mathbf{B}$, and the dramatic advantage of aRRR when the auxiliary signal $\mathbf{S}$ is strong. The latter point is critical, because molecular data typically have a large amount of structured variation that is driven by coordinated biological processes or other latent effects; it is common for such variation to be stronger than the signal of interest (i.e., $\mathbf{BY}$), yet it is not systematically adjusted for in practice.

For scenario (ii), we generate data $\{\mathbf{X}_j : 100 \times 100, \mathbf{Y}_j : 10 \times 100\}$ via $\mathbf{X}_j = (\mathbf{B}+\mathbf{B}_j)\mathbf{Y}+\mathbf{E}$ for two cohorts $j \in \{1,2\}$. Here, $\mathbf{B}_j$ are covariate effects specific to cohort $j$ and $\mathbf{B}$ are shared effects. The coefficient arrays are generated via $\mathbf{B} = a\mathbf{U}_B\mathbf{V}_B^T$, $\mathbf{B}_1 = b\mathbf{U}_{B_1}\mathbf{V}_{B_1}^T$, and $\mathbf{B}_2 = b\mathbf{U}_{B_2}\mathbf{V}_{B_2}^T$ where $\{\mathbf{U}_B, \mathbf{U}_{B_1}, \mathbf{U}_{B_2}\}$ are each $100 \times R_y$ and $\{\mathbf{V}, \mathbf{V}_{B_1}, \mathbf{V}_{B_2}\}$ are each $R_y \times 10$. The entries of $\{\mathbf{E}, \mathbf{Y}, \mathbf{U}_B, \mathbf{U}_{B_1}, \mathbf{U}_{B_2}, \mathbf{V}_B, \mathbf{V}_{B_1}, \mathbf{V}_{B_2}\}$ are each generated independently from a standard normal distribution. We consider $R_y = 1$ or 5, and three conditions with different signal strength for each term by adjusting $a$ and $b$: $a = 2$ and $b = 0.2$ ($||\mathbf{B}||/||\mathbf{B}_i|| = 10$), $a = b = 1$ (($||\mathbf{B}||/||\mathbf{B}_i|| = 1$), and $a = 0.2$ and $b = 2$ ($||\mathbf{B}||/||\mathbf{B}_i|| = 0.1$). For each set of conditions, we estimate $\mathbf{B}, \mathbf{B}_1$ and $\mathbf{B}_2$ for $J = 2$ via maRRR with no auxiliary terms $\mathbf{S}$, termed multi-cohort reduced rank regression (mRRR). Table 1B shows average relative MSEs of $\mathbf{B}$ and the $\mathbf{B}_i$'s for mRRR in comparison to two-stage approaches analogous to those described previously. The mRRR approach can effectively recover shared and cohort specific effects, with dramatic improvement over ad-hoc multi-step approaches.

Table 1: Relative MSE for scenarios assessing aRRR ($\mathbf{A}$) and mRRR ($\mathbf{B}$). Values that smaller than 0.01 are round to 0.01. The bold number represents the lowest value in a row.

| $\mathbf{A}$ | | aRRR | | SupSVD | | Two-stage LS | | Two-stage NN | |
|---|---|---|---|---|---|---|---|---|---|
| $\frac{||\mathbf{BY}||}{||\mathbf{S}||}$ | $R_y$ | $\mathbf{B}$ | $\mathbf{S}$ | $\mathbf{B}$ | $\mathbf{S}$ | $\mathbf{B}$ | $\mathbf{S}$ | $\mathbf{B}$ | $\mathbf{S}$ |
| 10 | 1 | **0.01** | 0.61 | 0.01 | **0.46** | 0.01 | 0.60 | 0.01 | 0.61 |
| 1 | 1 | **0.04** | 0.22 | 0.13 | **0.19** | 0.22 | 0.21 | 0.05 | 0.25 |
| 0.1 | 1 | **0.17** | **0.01** | 10.97 | 0.10 | 11.44 | 0.11 | 7.45 | 0.08 |
| 10 | 5 | **0.01** | 0.63 | 0.01 | **0.48** | 0.01 | 0.60 | 0.01 | 0.63 |
| 1 | 5 | **0.14** | 0.24 | 0.17 | **0.19** | 0.23 | 0.21 | 0.15 | 0.26 |
| 0.1 | 5 | **0.40** | **0.01** | 11.31 | 0.10 | 11.66 | 0.11 | 7.83 | 0.08 |

| $\mathbf{B}$ | | mRRR | | Two-stage LS | | Two-stage NN | |
|---|---|---|---|---|---|---|---|
| $\frac{||\mathbf{B}||}{||\mathbf{B}_i||}$ | $R_y$ | $\mathbf{B}$ | $\mathbf{B}_i$ | $\mathbf{B}$ | $\mathbf{B}_i$ | $\mathbf{B}$ | $\mathbf{B}_i$ |
| 10 | 1 | **0.01** | **0.11** | 0.01 | 0.77 | 0.01 | 0.42 |
| 1 | 1 | **0.01** | **0.01** | 0.75 | 0.60 | 0.67 | 0.52 |
| 0.1 | 1 | **0.07** | **0.01** | 80.57 | 0.55 | 76.17 | 0.52 |
| 10 | 5 | **0.01** | **0.28** | 0.01 | 0.57 | 0.01 | 0.50 |
| 1 | 5 | **0.08** | **0.08** | 0.49 | 0.54 | 0.45 | 0.50 |
| 0.1 | 5 | **0.49** | **0.01** | 54.79 | 0.56 | 52.41 | 0.54 |

## 6.2   Missing data imputation

In this simulation we assess the maRRR framework more broadly, with a focus on missing data imputation. Our general simulation procedure follows these steps: 1) complete data generation; 2) missingness assignment; 3) imputation analysis. In reality the true main signals may come from covariate effects or auxiliary structures and can be individual-level or shared across multiple cohorts. So we consider four fundamental scenarios: ($a$) large $\mathbf{B}$, main signals from one global auxiliary structure which is shared by all cohorts; ($b$) large $\mathbf{S}$, main signals

from one global covariate effect which is shared by all cohorts; ($c$) large $\mathbf{B}_i$, main signals from individual covariate effect of each cohort; ($d$) large $\mathbf{S}_i$, main signals from individual auxiliary structure in each individual cohort. In order to mimic the real situation, the number of samples and dimensions of the data is set to be the same as the TCGA data analyzed in Section 7. That is, $\mathbf{X}$ consists of 1000 features and 6581 samples from 30 study cohorts and $\mathbf{Y}$ consists of 50 predictors. Therefore, the ground truth can be written as $\mathbf{X}_j = \mathbf{B}\mathbf{Y}_j + \mathbf{B}_j\mathbf{Y}_j + \mathbf{S}_{shared,j} + \mathbf{S}_j + \mathbf{E}_j, j = 1, ..., 30$. In each simulation, the standard deviation for the main signals is set to be $\sqrt{10}$ while that of the remaining signals and random errors are set to be 1. The complete data generation process is described in Appendix D.1.

To mimic the various types of missingness that are encountered in reality, we conduct simulations in which four types of missingness are considered: (i) missing entries, (ii) missing columns, (iii) missing rows, and (iv) a balanced mix of these three types of missingness as the average of the results of those three types. Missingness is set to be 5% of the original data for the assumption that adequate information is provided for revealing global structures. All missing indices are randomly selected. Denote $\widetilde{\mathbf{X}}.$ as the estimate for true observation $\mathbf{X}.$, based on non-missing entries. We define the relative squared error (RSE) for missing data imputation as $RSE = \frac{\sum_{(m,n)\in M}(\mathbf{X}.[(m,n)]-\widetilde{\mathbf{X}}.[(m,n)])^2}{\sum_{(m,n)\in M}(\mathbf{X}.[(m,n)])^2}$.

We compare our method (maRRR with true 31 modules) with the following approaches: (1) BIDIFAC+ with 31 modules, i.e. only auxiliary variation structure $\mathbf{S}$; (2) mRRR with 31 modules, i.e. only covariate-related structure $\mathbf{BY}$; (3) aRRR with only one module for all cancer types' cohorts together; (4) aRRR separately on each cancer type's cohort; (5) nuclear norm regression (without $\mathbf{S}$) of $\mathbf{X}.$ on $\mathbf{Y}.$, i.e. mRRR with one all-shared module; (6) nuclear norm regression (without $\mathbf{S}$) of $\mathbf{X}_j$ on $\mathbf{Y}_j$ separate for each cancer type, i.e. mRRR with 30 individual modules; (7) nuclear norm approximation (without $\mathbf{BY}$) for all cancer types together (8) nuclear norm approximation (without $\mathbf{BY}$) for each cancer types separately.

Based on the simulation results shown in Table 2, in terms of the average performance, our proposed method maRRR has the lowest RSE. In particular, maRRR has a very close RSE to the best result in the case of missing columns or rows under large individual covariate signals, and it performs the best in all other cases. BIDIFAC+ performs slightly worse than our proposed method while there are only missing entries or rows, but it cannot utilize any covariate information to impute in the case of missing columns. The models only considering covariate effects (mRRR and nuclear norm regression) cannot predict accurately when the auxiliary variation is large, no matter global or individual. On the contrary, the models without considering covariates (BIDIFAC+, nuclear norm approximation) can perform reasonably well in the cases of missing entries or rows since the variation from covariates may be counted into that of auxiliary structures. The special case of our proposed method, aRRR (for one cohort only), though worse than maRRR, has lower MSE than many other existing methods.

## 6.3 Computation

Our proposed method is computationally efficient. For a matrix of size $1000 \times 6581$, the average computational cost per epoch is 45 seconds for Algorithm 1 and 50 seconds for Algorithm 2. A comprehensive comparison of computation times for all methods utilized in this study is provided in Appendix D.2.

11

Table 2: Imputation relative squared error(RSE) under different methods and different types of missingness, simulated data is set to be large at only one type of modules. Missingness is set to be 5% of the original **X**. The number of epochs for each method is set as 30. Each result is a mean of 10 replications. The standard error is less than 0.01 for all of the means shown. The bold number represent the lowest value in a column.

| large_B | Method | missing entries | missing columns | missing rows | mean |
|---|---|---|---|---|---|
| | maRRR | **0.082** | **0.228** | **0.216** | **0.175** |
| | BIDIFAC+ | 0.085 | 1 | 0.231 | 0.439 |
| | mRRR | 0.202 | 0.241 | 0.285 | 0.243 |
| | aRRR, one all-shared | 0.125 | 0.288 | 0.227 | 0.213 |
| | aRRR, 30 separate | 0.093 | 0.287 | 1.014 | 0.465 |
| | NN reg, one all-shared | 0.283 | 0.287 | 0.288 | 0.286 |
| | NN reg, 30 separate | 0.212 | 0.255 | 1 | 0.489 |
| | NN approx, one all-shared | 0.127 | 1 | 0.229 | 0.452 |
| | NN approx, 30 separate | 0.096 | 1 | 1 | 0.699 |
| large_S | Method | missing entries | missing columns | missing rows | mean |
| | maRRR | **0.082** | **0.877** | **0.218** | **0.392** |
| | BIDIFAC+ | 0.085 | 1 | 0.225 | 0.437 |
| | mRRR | 0.759 | 1.066 | 0.927 | 0.917 |
| | aRRR, one all-shared | 0.125 | 0.929 | 0.228 | 0.427 |
| | aRRR, 30 separate | 0.093 | 0.884 | 1.004 | 0.66 |
| | NN reg, one all-shared | 0.928 | 0.931 | 0.933 | 0.93 |
| | NN reg, 30 separate | 0.783 | 1.072 | 1 | 0.952 |
| | NN approx, one all-shared | 0.127 | 1 | 0.23 | 0.452 |
| | NN approx, 30 separate | 0.096 | 1 | 1 | 0.699 |
| large_Bi | Method | missing entries | missing columns | missing rows | mean |
| | maRRR | **0.083** | 0.262 | 0.901 | **0.415** |
| | BIDIFAC+ | 0.086 | 1 | **0.867** | 0.651 |
| | mRRR | 0.204 | **0.246** | 0.928 | 0.459 |
| | aRRR, one all-shared | 0.149 | 0.913 | 0.902 | 0.655 |
| | aRRR, 30 separate | 0.093 | 0.287 | 1.013 | 0.464 |
| | NN reg, one all-shared | 0.896 | 0.899 | 0.964 | 0.92 |
| | NN reg, 30 separate | 0.212 | 0.252 | 1 | 0.488 |
| | NN approx, one all-shared | 0.151 | 1 | 0.907 | 0.686 |
| | NN approx, 30 separate | 0.096 | 1 | 1 | 0.699 |
| large_Si | Method | missing entries | missing columns | missing rows | mean |
| | maRRR | **0.083** | **0.873** | **0.861** | **0.606** |
| | BIDIFAC+ | 0.086 | 1 | 0.866 | 0.651 |
| | mRRR | 0.76 | 1.066 | 0.928 | 0.918 |
| | aRRR, one all-shared | 0.148 | 0.93 | 0.906 | 0.661 |
| | aRRR, 30 separate | 0.093 | 0.885 | 1.003 | 0.661 |
| | NN reg, one all-shared | 0.929 | 0.931 | 0.935 | 0.932 |
| | NN reg, 30 separate | 0.784 | 1.072 | 1 | 0.952 |
| | NN approx, one all-shared | 0.15 | 1 | 0.91 | 0.687 |
| | NN approx, 30 separate | 0.096 | 1 | 1 | 0.699 |

# 7 Real Data Analysis

## 7.1 Data description

We consider data from the Cancer Genome Atlas (TCGA) Pan-Cancer Project (Hoadley et al., 2018). TCGA is an NIH-sponsored initiative to molecularly characterize cancer tissue samples obtained from hundreds of sites worldwide. We used data for 6581 tumor samples from distinct individuals representing 30 different cancer types (i.e., 30 cohorts). The number of samples per cancer type ranges from 57 samples for uterine carcinosarcoma (UCS) to 976 for breast carcinoma (BRCA). As outcomes, we consider gene expression data obtained from Illumina RNASeq platforms and normalized as described in Hoadley et al. (2018). We filter to the 1000 genes that have the highest standard deviation, yielding $\mathbf{X}. : 1000 \times 6581$. We filter to the 50 somatic most common somatic mutations (1=mutated and 0=not mutated) as covariates $\mathbf{Y}. : 50 \times 6581$. Data are standardized and orthogonalized as in Appendix C. Further details on the data are available in Appendix E.1 and E.2.

## 7.2 Decomposition results

We first apply the optimization with dynamic modules for BIDIFAC+ (Lock et al., 2022), to uncover 50 low-rank modules in $\mathbf{X}.$. This stepwise procedure iteratively determines modules of shared structure without consideration of any covariates (i.e., $\mathbf{C}_S$). Fifty was chosen as the upper bound because the variance explained by more modules than 50 was relatively inconsequential. To distinguish how much variance is attributed to mutation effects, we set covariate-related module indicators $\mathbf{C}_Y$ equal to $\mathbf{C}_S$. We then apply maRRR to these 50 modules to infer both mutation-driven and auxiliary structured variation, ($K = L = 50$) with penalties determined as in Section 4.

We order the modules by total variance explained by descending by maRRR estimates, i.e. $||\widehat{\mathbf{B}}_i\mathbf{Y}.^{(i)} + \widehat{\mathbf{S}}.^{(i)}||_2^2, i = 1, ..., 50$. The ordered result is shown in Table 3. The top 3 modules by total variance explained are those with one or two cancers: BRCA, THCA and a combination of GBM and LGG (both neurological cancers), respectively. In general, auxiliary structures explained more variation than mutation-related structures, but their relative contribution varied widely across modules. For example, Modules 6 and 7 have fairly comparable amount of variation explained by both mutation-related and -unrelated parts. Other modules have negligible mutation-driven variation, such as Module 12 which is shared by all but LAML. The large amount of low-rank variation unrelated to covariates demonstrates the importance of accounting for this auxiliary structure. Moreover, the large amount of covariate-related and -unrelated variation that is specific to one or a small number of cancer types demonstrates the importance of accounting for individual and partially-shared structures.

We present a comparison of the BIDIFAC+ and mRRR estimates with those of our proposed maRRR in Table 3. The total signal detected by maRRR closely aligns with that of BIDIFAC+, illustrating maRRR's ability to discern covariate effects while accounting for similar total variance. The covariate effects identified by mRRR resemble those by maRRR, particularly when the sample size is large. However, mRRR tends to estimate larger covariate effects, especially for smaller sample sizes. This makes sense because the maRRR estimates are less prone to over-fitting by adjusting for unrelated structure. As the sample size decreases, the relative square errors (RSE) between estimates from the two methods increase. However, both methods generally agree when the mutation effects are negligible ($\leq 10000$). For instance, both mRRR and maRRR reveal virtually no global mutation effects for Module 12 (even though the RSE is large due to the relative standardization).

Principal components plots of the first three modules (BRCA, THCA, GBM and LGG) are shown in Figure 2. Figure 2A displays mutation-related variation for the BRCA module,

Table 3: Cancer types and sources for the first 15 modules, ordered by variation explained by maRRR. Variance of $\mathbf{S}^{(i)}_{\cdot}/\mathbf{B}_i\mathbf{Y}^{(i)}_{\cdot}/\mathbf{B}_i\mathbf{Y}^{(i)}_{\cdot}$/signal refers to total variance explained by $\mathbf{S}^{(i)}_{\cdot}/\mathbf{B}_i\mathbf{Y}^{(i)}_{\cdot}/\mathbf{B}_i\mathbf{Y}^{(i)}_{\cdot}+\mathbf{S}^{(i)}_{\cdot}$ estimated by maRRR and mRRR. RSE of BIDIFAC+ refers to relative square difference between signal ($\mathbf{B}_i\mathbf{Y}^{(i)}_{\cdot}+\mathbf{S}^{(i)}_{\cdot}$) estimated by maRRR and BIDIFAC+. Sample size refers to the number of samples used in the current module. All the number smaller than 0.01 are round to 0.01.

For reference, the RSE for mRRR and maRRR is defined as $\frac{\|\mathbf{S}^{(i)}_{\cdot,BIDIFAC+}-\mathbf{S}^{(i)}_{\cdot,maRRR}\|^2_F}{\|\mathbf{S}^{(i)}_{\cdot,BIDIFAC+}+\mathbf{S}^{(i)}_{\cdot,maRRR}\|^2_F}$ for any $i=1,2,...,50$.

the RSE for mRRR and maRRR is defined as $\frac{\|\mathbf{B}_{i,mRRR}\mathbf{Y}^{(i)}_{\cdot}-\mathbf{B}_{i,maRRR}\mathbf{Y}^{(i)}_{\cdot}\|^2_F}{\|\mathbf{B}_{i,mRRR}\mathbf{Y}^{(i)}_{\cdot}+\mathbf{B}_{i,maRRR}\mathbf{Y}^{(i)}_{\cdot}\|^2_F}$ since $\mathbf{Y}^{(i)}_{\cdot}$ is semi-orthogonal;

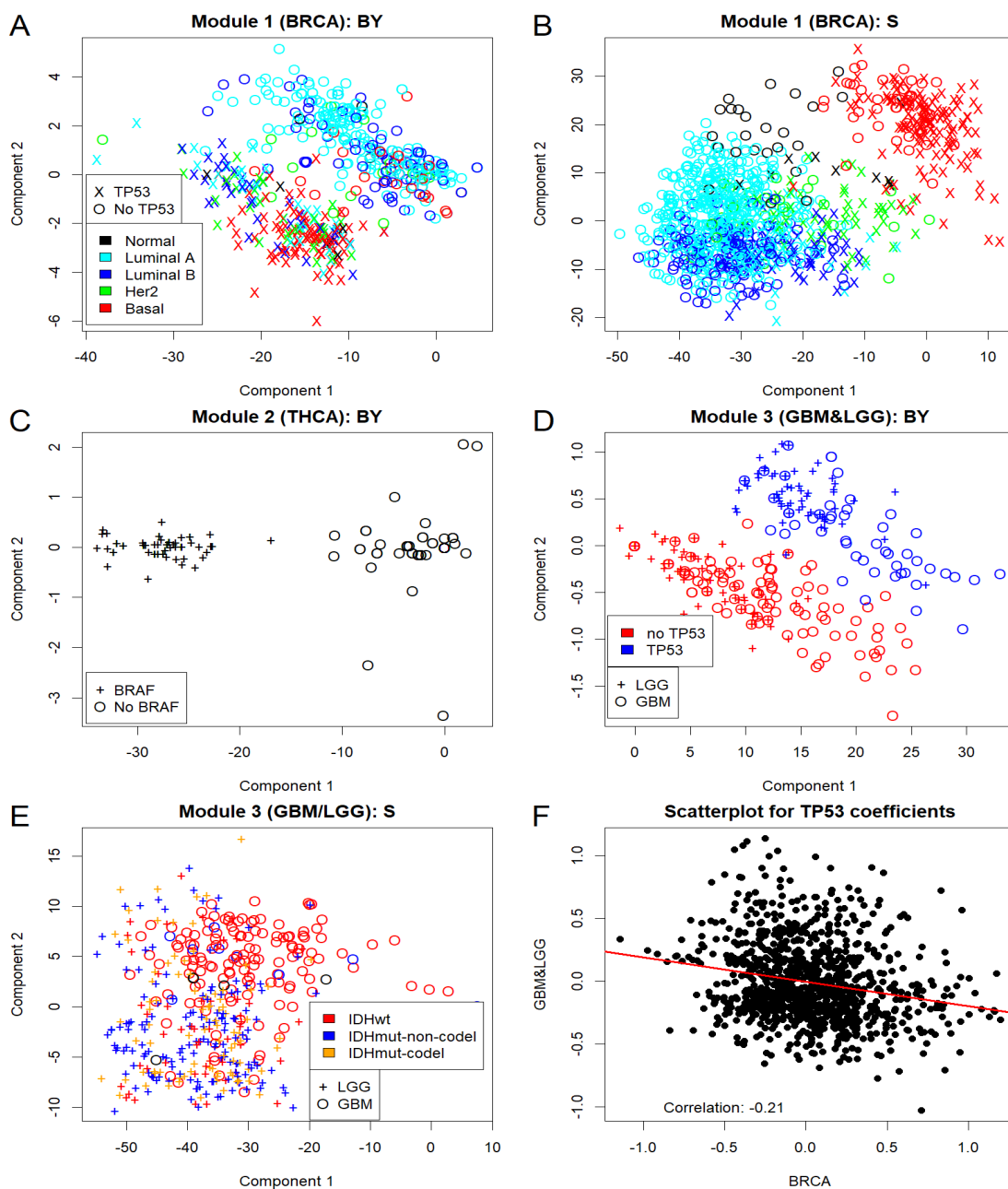| Module ($i$) | Sample size | Variance of $\mathbf{B}_i\mathbf{Y}^{(i)}_{\cdot}$ | RSE of mRRR | Variance of $\mathbf{S}^{(i)}_{\cdot}$ | Variance of signal | RSE of BIDIFAC+ | Cancer types |
|---|---|---|---|---|---|---|---|
| 1 | 976 | 129652.09 | 0.22 | 1085292.83 | 1524117.23 | 0.01 | BRCA |
| 2 | 400 | 167384.83 | 0.15 | 551269.19 | 992693.82 | 0.01 | THCA |
| 3 | 433 | 67242.88 | 0.20 | 645839.92 | 975136.09 | 0.01 | GBM, LGG |
| 4 | 331 | 49.52 | 0.97 | 730462.99 | 734935.73 | 0.01 | PRAD |
| 5 | 195 | 7231.61 | 0.68 | 633034.84 | 724960.62 | 0.01 | LIHC |
| 6 | 1029 | 137934.04 | 0.08 | 316746.84 | 656053.3 | 0.03 | BLCA, CESC, ESCA, HNSC, KICH, LUSC |
| 7 | 820 | 137765.92 | 0.14 | 298195.14 | 629007.34 | 0.02 | COAD, ESCA, PAAD, READ, STAD |
| 8 | 179 | 20.61 | 1.00 | 396671.22 | 396882.04 | 0.01 | PCPG |
| 9 | 170 | 53.04 | 0.98 | 347389.29 | 348257.22 | 0.01 | LAML |
| 10 | 576 | 13530.5 | 0.33 | 231107.59 | 316968.61 | 0.04 | KIRC, KIRP |
| 11 | 422 | 30138.96 | 0.21 | 192612.5 | 300435.91 | 0.03 | SKCM, UVM |
| 12 | 6411 | 0.14 | 0.74 | 264645.95 | 264654.13 | 0.23 | All but *LAML* |
| 13 | 211 | 63836.59 | 0.16 | 87452.29 | 251028.58 | 0.01 | COAD, READ |
| 14 | 149 | 495.86 | 0.87 | 237614.82 | 250811.31 | 0.01 | TGCT |
| 15 | 119 | 97.81 | 0.97 | 242978.25 | 243471.26 | 0.01 | THYM |

14

Figure 2: **A**: scores for the first two principle components of covariate-related variation (**BY**) from Module 1 (BRCA); **B**: scores for the first two principle components of covariate-unrelated auxiliary variation (**S**) from Module 1 (BRCA), with symbols and colors showing TP53 mutation and 5 subtypes of BRCA respectively. **C**: scores for the first two principle components of covariate-related variation (**BY**) from Module 2 (THCA), with symbols showing BRAF mutation. Plot **A**&**B** contain 976 samples in the BRCA cohort respectively; Plot **C** contains 400 samples in the BRCA cohort respectively. **D**: scores for the first two principle components of covariate-related variation (**BY**) from Module 3 (GBM&LGG); **E**: scores for the first two principle components of covariate-unrelated auxiliary variation (**S**) from Module 3 (GBM&LGG), with symbols and colors showing TP53 mutation and 3 IDH/codel subtypes of GBM&LGG respectively. **F**: scatterplot for the regression coefficients for TP53 in module 1 and module 3. Each point in Plot **D**&**E** represents one sample (150 for GBM and 283 for LGG); each point in Plot **F** represents one gene (1000 in total).

and samples are distinguished by whether they have a mutation in the TP53 gene or not. This makes sense, as TP53 is known to play a critical role in genomic activity for cancer and it is the most frequently mutated gene in breast cancer (Olivier et al., 2010). From Figure 2B,we see that the mutation-unrelated structure is driven by the 5 intrinsic BRCA subtypes (Cancer Genome Atlas Research Network, 2012): Normal-like, Luminal A (LumA), Luminal B (LumB), HER2-enriched (HER2), and Basal-enriched (Basal) tumors. This makes sense, as the BRCA subtypes are known to be genomically distinct, but (as is apparent) do not have a direct correspondence to TP53 or other common mutations. From Figure 2C, we observe that BRAF and non-BRAF groups are elegantly separated on the first principle component of mutation-driven variation for THCA, which explains much more variation than other components. This concurs with prior research indicating that the BRAF mutation defines a unique genomic and clinical subgroup within THCA patients (Dolezal et al., 2021).

Figure 2D&E reveal substantial variation in both the GBM and LGG samples, as evidenced by the spread and intermingling of the "+" and "o" symbols. This observation indicates that the top components identified by Module 3 explain substantial variation in both the GBM and LGG cancer cohorts, suggesting that it is indeed shared by the two cancer types. The TP53 mutation drives separation of the mutation-driven structure, which makes sense as TP53 status is closely related to GBM and LGG aggressiveness (Ham et al., 2019). This was also observed in BRCA, however, from Figure 2F we see that the TP53 regression coefficients from module 3 (GBM&LGG) has little correlation to those from module 1 (BRCA). This demonstrates that while TP53 is an important somatic mutation related to different types of cancer, its effect on gene expression can differ dramatically depending on the cancer type. In general, we find that the same somatic mutation plays different roles for different cohorts and modules. This embodies the necessity of the flexible modeling for different (combinations of) cohorts. In fact, one interesting finding of our analysis is that the effect of somatic mutations on gene expression are almost not entirely shared across different types of cancer. For example, for the module that is shared across almost all cancers (module 12) the mutation-driven component is negligible. This observation is further supported by our analysis in Appendix E.4 and Section 7.3.

## 7.3   Missing data imputation

Similar to Section 6.2, we compare our proposed maRRR with other relevant methods under four types of missingness for these data. Beyond the aforementioned eight methods in Section 6.2, we added (9) linear least squares regression to predict $\mathbf{X}$ from $\mathbf{Y}$ for all cancer types together and (10) linear least-squares regression for each cancer types separately. Note that maRRR, BIDIFAC+ and mRRR are based on the detected 50 modules. Results are provided in Table 4. In the scenario of missing entries, both maRRR and BIDIFAC+ have the lowest RSE, with similar values; both methods allow for an efficient decomposition of joint and individual structures. In the case of missing columns (some samples' entire outcomes are missing), the methods that do not consider mutations and only consider $\mathbf{S}$ (BIDIFAC+, NN approx) have no predictive power, which is expected because the mutation data is needed to inform predictions if no gene expression is available for a sample. Here the methods that consider both $\mathbf{B}$ and $\mathbf{S}$ (maRRR, aRRR) are suboptimal compared to those methods that only consider mutation-driven variation ($\mathbf{BY}$), indicating that for these data allowing for auxiliary variation does not improve column-wise predictions. Moreover, methods that allow for separate mutation-driven structure across the cohorts perform substantially better, which is consistent with the fact that there were more individual modules in our analysis and mutation-driven variation was generally not shared (Table 3). In the event of missing rows (each cohort misses random features), methods that consider individual structure only do not perform well, as

Table 4: Imputation relative squared error(RSE) under different methods and different types of missingness. Missingness is set to be 5% of the original **X**. "one all shared" means data for 30 groups are stacked together to form one matrix to analyze; "30 separate" means each group has its only model. "Missing entries" refers to missingness is entrywise; "missing columns" means some samples' entire observation are missing; "missing rows" means each group has several features entirely missing. "N/A" means some specific method is not applicable. The bold number represents the lowest value in a column.

| Methods | Missing entries | Missing columns | Missing rows | Average |
|---|---|---|---|---|
| maRRR | **0.233** | 0.813 | 0.600 | **0.548** |
| BIDIFAC+ | **0.233** | 0.999 | 0.613 | 0.615 |
| mRRR | 0.603 | **0.711** | 0.998 | 0.770 |
| aRRR, one all-shared | 0.261 | 0.930 | 0.487 | 0.559 |
| aRRR, 30 separate | 0.376 | 0.780 | 1.001 | 0.719 |
| LS reg, one all-shared | 0.908 | 0.906 | 0.899 | 0.904 |
| LS reg, 30 separate | 0.560 | N/A | N/A | N/A |
| NN reg, one all-shared | 0.912 | 0.913 | 1.032 | 0.953 |
| NN reg, 30 separate | 0.599 | 0.727 | 1.000 | 0.775 |
| NN approx, one all-shared | 0.273 | 1.000 | **0.454** | 0.576 |
| NN approx, 30 separate | 0.252 | 1.000 | 1.009 | 0.754 |

they cannot leverage shared structure when a gene is entirely missing within a cohort. On the contrary, methods with only one all-shared module (NN approx and aRRR) perform well. Here, maRRR also performs reasonable well, as including several individual modules does not limit its performance. In this case methods considering covariate effects only (mRRR, NN reg) do not perform well, as they tend toward estimates of zero (i.e., no predictions) to minimize squared error loss. NN approx is slightly better than aRRR., perhaps because it does not consider covariate-driven variation.

Under the circumstances of a balanced mix of missingness for different conditions, maRRR has the best average recovering ability. This is largely because it is the most robust and flexible. Other comparable methods (BIDIFAC+, aRRR, NN approx) will have limited peformnance for at least one form of missingness. In reality, maRRR will be the most suitable for imputation since missingness is unpredictable and complex.

# 8    Discussion

Two strengths of our proposed maRRR approach are its flexibility and versatility. It is flexible because it accounts for various types of signals - covariate-driven, shared or unshared - without prior assumptions on the size or rank of these signals. It is versatile because it is capable of performing many tasks at once: e.g., dimension reduction, prediction and missing data imputation. These advantages are well-illustrated by our pan-cancer application, in which adequate amounts of variation are explained by different components and the patterns detected are both insightful and consistent with existing scientific research on cancer.

We focus on multi-cohort integration rather than multi-view (data on the same subjects from different sources) integration, in part because shared or unshared covariate effects are straightforward to interpret across multiple cohorts. But one can still argue that in a multi-view (e.g., multi-omics) context each sample will have intrinsic underlying signals that will affect variables from different sources. Without loss of generality, this method can be adapted to analyze multi-view data as well. This is achieved by simply switching the way we integrate

matrices: horizontally across shared rows or vertically across shared columns. A promising future direction is to extend maRRR to the bidimensional integration context, where the data are both multi-cohort and multi-view. Another direction of future work is alternative empirical approaches to determine the module indicator matrices $\mathbf{C}_Y$ and $\mathbf{C}_S$, such as via an iterative stepwise selection procedure extending that in (Lock et al., 2022). While we have fixed the selection of penalty parameters $\lambda_B^{(k)}, \lambda_S^{(l)}$ by employing random matrix theory, the parameters or the ranks of the underlying structures may be estimated empirically by a cross-validation procedure combined with a grid search.

Further theoretical developments are also a pertinent future direction, such as proving the convergence of our optimization algorithms to a global optimum and establishing sufficient conditions for the uniqueness of the solution. There is empirical evidence for both conjectures, as we find that the converged solution is the same with different initializations and for the two optimization algorithms considered. Moreover, Theorem 1 of Lock et al. (2022) provides sufficient conditions for conditional uniqueness of the $\{\mathbf{S}^{(l)}\}_{l=1}^L$ given $\{\mathbf{B}_k\}_{k=1}^K$, and vice-versa.

## Data Availability Statement

The data that support the findings in this paper are provided via this RData file. The user-friendly R package maRRR at `https://github.com/JiuzhouW/maRRR` performs all functions described herein, such as fitting models by the two algorithms in Section 5.1, imputing missing values as in Section 6.2, generating penalties as in Section 4, and generating data as in Appendix D.1. For real data analysis, we provide all the model estimates as Rdata file with detailed notation explanations and heatmaps for all module estimates in an online file.

## Acknowledgements

# Appendix A   Additional methodological details

## A.1   Notation details

Detailed explanations for our proposed model 1 in the main manuscript are listed in Table 5.

Table 5: Notation for the proposed model.

| | | |
|---|---|---|
| $J$ | Observed | Number of cohorts |
| $K$ | Pre-specified | Number of covariate effects |
| $L$ | Pre-specified | Number of auxiliary structures |
| $\mathbf{X}_j$ | Observed | Outcome matrix for $j$th cohort |
| $\mathbf{Y}_j$ | Observed | Covariate matrix for $j$th cohort |
| $\mathbf{Y}_{\cdot}^{(k)}$ | Pre-specified | Design matrix for $k$th covariate effect concatenated from all $J$ cohorts |
| $\mathbf{B}_k$ | Estimated | coefficients for $k$th covariate effect |
| $\mathbf{S}_{\cdot}^{(l)}$ | Estimated | $l$th auxiliary structure concatenated from all $J$ cohorts |
| $\mathbf{E}_j$ | Estimated | Random error matrix for the $j$th cohort |
| $\mathbf{C}_Y$ | Pre-specified | Binary indicator matrix where its $[j,k]$th entry |
| | | determines whether $j$th cohort is considered in $k$th covariate effect |
| $\mathbf{C}_S$ | Pre-specified | Binary indicator matrix where its $[j,l]$th entry |
| | | determines whether $j$th cohort is considered in $l$th auxiliary structure |
| $\mathbf{U}_S^{(l)}$ | Estimated | Loading matrix of $l$th auxiliary structure |
| $\mathbf{V}_{Sj}^{(l)}$ | Estimated | Score matrix of $l$th auxiliary structure for $j$th cohort |
| $\mathbf{U}_B^{(k)}$ | Estimated | Loading matrix of $k$th covariate coefficients |
| $\mathbf{V}_B^{(k)}$ | Estimated | Score matrix of $k$th covariate coefficients |

## A.2   Construction of module indicator matrices

The construction of module indicator matrices $\mathbf{C}_Y$ and $\mathbf{C}_S$ can be accomplished in various ways depending on the specific context and available prior knowledge. If there exists prior knowledge indicating shared effects among certain cohorts, it would be straightforward to define modules accordingly. For example, defining a global module and individual modules for each cohort, as illustrated in Section 6.2. In absence of such information, there are other practical methods that can be applied.

In scenarios where the number of cohorts is small, one can begin by enumerating all possible combinations of cohorts in $\mathbf{C}_Y$ and $\mathbf{C}_S$. As the objective function encourages rank sparsity, modules with no true shared structure may be estimated as zero (i.e., no structure) even if they are included in the algorithm. Moreover, after obtaining initial estimates, modules that explain a relatively higher amount of variance can be retained to form updated $\mathbf{C}_Y$ and $\mathbf{C}_S$.

When dealing with a large number of cohorts, an alternative approach would be to adopt a data-driven strategy like the "Optimization algorithm: dynamic modules" section from Lock et al. (2022). This method is designed to select $\mathbf{C}_S$ based on the amount of variance explained, after which $\mathbf{C}_Y$ can be set to match $\mathbf{C}_S$, thus partitioning the variance related to covariates. This is illustrated in Section 7.2. Optionally as a second step, one could keep the modules that explain a high amount of variance to reformulate $\mathbf{C}_Y$ and $\mathbf{C}_S$. This can help identify the most significant components in the covariate-related effects and auxiliary structures.

# Appendix B   Proofs

## B.1   Proof of Theorem 1

*Proof.* Consider the following lemma, the proof of which is provided in Mazumder et al. (2010):

**Lemma 2.** *(Mazumder et al., 2010) For any matrix $\mathbf{Z} : m \times n$ with $rank(\mathbf{Z}) = k$, $\forall r \geq k$, the following holds:*

$$||\mathbf{Z}||_* = \min_{\substack{\mathbf{U},\mathbf{V}: \\ \mathbf{Z}=\mathbf{U}_{m\times r}\mathbf{V}_{n\times r}^T}} \frac{1}{2}(||\mathbf{U}||_F^2 + ||\mathbf{V}||_F^2)$$

Applying Lemma 2 to each $\mathbf{B}_k, k = 1, ..., K$ and $\mathbf{S}_{\cdot}^{(l)}, l = 1, ..., L$:

$$\min_{\{\mathbf{B}_k\}_{k=1}^K,\{\mathbf{S}_{\cdot}^{(l)}\}_{l=1}^L} \{\frac{1}{2}||\mathbf{X}_{\cdot} - \sum_{k=1}^K \mathbf{B}_k\mathbf{Y}_{\cdot}^{(k)} - \sum_{l=1}^L \mathbf{S}_{\cdot}^{(l)}||_F^2 + \sum_{k=1}^K \lambda_B^{(k)}||\mathbf{B}_k||_* + \sum_{l=1}^L \lambda_S^{(l)}||\mathbf{S}_{\cdot}^{(l)}||_*\}$$

$$= \min_{\{\mathbf{B}_k\}_{k=1}^K,\{\mathbf{S}_{\cdot}^{(l)}\}_{l=1}^L} \{\frac{1}{2}||\mathbf{X}_{\cdot} - \sum_{k=1}^K \mathbf{B}_k\mathbf{Y}_{\cdot}^{(k)} - \sum_{l=1}^L \mathbf{S}_{\cdot}^{(l)}||_F^2 +$$

$$\sum_{k=1}^K \lambda_B^{(k)} \min_{\mathbf{U}_B^{(k)},\mathbf{V}_B^{(k)}:\mathbf{B}_k=\mathbf{U}_B^{(k)}\mathbf{V}_B^{(k)T}} \frac{1}{2}(||\mathbf{U}_B^{(k)}||_F^2 + ||\mathbf{V}_B^{(k)}||_F^2) +$$

$$\sum_{l=1}^L \lambda_S^{(l)} \min_{\mathbf{U}_S^{(l)},\mathbf{V}_S^{(l)}:\mathbf{S}_{\cdot}^{(l)}=\mathbf{U}_S^{(l)}\mathbf{V}_S^{(l)T}} \frac{1}{2}(||\mathbf{U}_S^{(l)}||_F^2 + ||\mathbf{V}_S^{(l)}||_F^2)\}$$

$$= \min_{\substack{\{\mathbf{U}_B^{(k)},\mathbf{V}_B^{(k)}:\mathbf{B}_k=\mathbf{U}_B^{(k)}\mathbf{V}_B^{(k)T}\}_{k=1}^K, \\ \{\mathbf{U}_S^{(l)},\mathbf{V}_S^{(l)}:\mathbf{S}_{\cdot}^{(l)}=\mathbf{U}_S^{(l)}\mathbf{V}_S^{(l)T}\}_{l=1}^L}} \frac{1}{2}\{||\mathbf{X}_{\cdot} - \sum_{k=1}^K \mathbf{U}_B^{(k)}\mathbf{V}_B^{(k)T}\mathbf{Y}_{\cdot}^{(k)} - \sum_{l=1}^L \mathbf{U}_S^{(l)}\mathbf{V}_S^{(l)T}||_F^2 +$$

$$\sum_{k=1}^K \lambda_B^{(k)} \min_{\mathbf{U}_B^{(k)},\mathbf{V}_B^{(k)}:\mathbf{B}_k=\mathbf{U}_B^{(k)}\mathbf{V}_B^{(k)T}}(||\mathbf{U}_B^{(k)}||_F^2 + ||\mathbf{V}_B^{(k)}||_F^2) +$$

$$\sum_{l=1}^L \lambda_S^{(l)} \min_{\mathbf{U}_S^{(l)},\mathbf{V}_S^{(l)}:\mathbf{S}_{\cdot}^{(l)}=\mathbf{U}_S^{(l)}\mathbf{V}_S^{(l)T}}(||\mathbf{U}_S^{(l)}||_F^2 + ||\mathbf{V}_S^{(l)}||_F^2)\}$$

$$= \min_{\{\mathbf{U}_B^{(k)},\mathbf{V}_B^{(k)}\}_{k=1}^K,\{\mathbf{U}_S^{(l)},\mathbf{V}_S^{(l)}\}_{l=1}^L} \frac{1}{2}\{||\mathbf{X}_{\cdot} - \sum_{k=1}^K \mathbf{U}_B^{(k)}\mathbf{V}_B^{(k)T}\mathbf{Y}_{\cdot}^{(k)} - \sum_{l=1}^L \mathbf{U}_S^{(l)}\mathbf{V}_S^{(l)T}||_F^2 +$$

$$\sum_{k=1}^K \lambda_B^{(k)}(||\mathbf{U}_B^{(k)}||_F^2 + ||\mathbf{V}_B^{(k)}||_F^2) + \sum_{l=1}^L \lambda_S^{(l)}(||\mathbf{U}_S^{(l)}||_F^2 + ||\mathbf{V}_S^{(l)}||_F^2).\}$$

∎

## B.2   Proof of Proposition 1

*Proof.* Assume a violation of condition 1, wherein $\lambda_B^{(k)} \geq \frac{1}{c_y}\sum_{i\in\mathcal{I}_k}\lambda_B^{(i)}$. Let $\widehat{\mathbf{B}}_k\mathbf{Y}_{\cdot}^{(k)} = \sum_{i\in\mathcal{I}_k}\widehat{\mathbf{B}}_i'\mathbf{Y}_{\cdot}^{(i)}$, where $\widehat{\mathbf{B}}_i'\mathbf{Y}_{\cdot}^{(i)}$ contains the blocks of $\widehat{\mathbf{B}}_k\mathbf{Y}_{\cdot}^{(k)}$ corresponding to $\mathbf{C}_Y[\cdot,i]$ and $\mathbf{0}$ otherwise. The choice of $\widehat{\mathbf{B}}_i'$ is unique that $\widehat{\mathbf{B}}_i' = \frac{1}{c_y}\widehat{\mathbf{B}}_k$. For all $i \in \mathcal{I}_k$, we have $||\widehat{\mathbf{B}}_i'\mathbf{Y}_{\cdot}^{(i)}||_* = ||\frac{1}{c_y}\widehat{\mathbf{B}}_k\mathbf{Y}_{\cdot}^{(i)}||_* \leq ||\frac{1}{c_y}\widehat{\mathbf{B}}_k\mathbf{Y}_{\cdot}^{(k)}||_* \leq ||\widehat{\mathbf{B}}_k\mathbf{Y}_{\cdot}^{(k)}||_*$ since $c_y \geq 1$. Consider a minimizer $\{\widetilde{\mathbf{B}}_k\}_{k=1}^K, \{\widetilde{\mathbf{S}}_{\cdot}^{(l)}\}_{l=1}^L$,

where $\{\widetilde{\mathbf{S}}_{\cdot}^{(l)}\}_{l=1}^{L} = \{\widehat{\mathbf{S}}_{\cdot}^{(l)}\}_{l=1}^{L}$, $\widetilde{\mathbf{B}}_k = \mathbf{0}$, $\widetilde{\mathbf{B}}_i = \widehat{\mathbf{B}}_i + \widehat{\mathbf{B}}_i'$, $\forall i \in \mathcal{I}_k$, and all other $\mathbf{B}$ estimates are equal. Then, by the triangle inequality,

$$
\begin{aligned}
f(\{\widehat{\mathbf{B}}_k\}_{k=1}^{K}, \{\widehat{\mathbf{S}}_{\cdot}^{(l)}\}_{l=1}^{L}) - f(\{\widetilde{\mathbf{B}}_k\}_{k=1}^{K}, \{\widetilde{\mathbf{S}}_{\cdot}^{(l)}\}_{l=1}^{L}) &= \lambda_B^{(k)}||\widehat{\mathbf{B}}_k||_* + \sum_{i \in \mathcal{I}_k} \lambda_B^{(i)}||\widehat{\mathbf{B}}_i||_* - \sum_{i \in \mathcal{I}_k} \lambda_B^{(i)}||\widehat{\mathbf{B}}_i + \widehat{\mathbf{B}}_i'||_* \\
&\geq \lambda_B^{(k)}||\widehat{\mathbf{B}}_k||_* + \sum_{i \in \mathcal{I}_k} \lambda_B^{(i)}||\widehat{\mathbf{B}}_i||_* - \sum_{i \in \mathcal{I}_k} \lambda_B^{(i)}(||\widehat{\mathbf{B}}_i||_* + ||\widehat{\mathbf{B}}_i'||_*) \\
&= \lambda_B^{(k)}||\widehat{\mathbf{B}}_k||_* - \sum_{i \in \mathcal{I}_k} \lambda_B^{(i)}||\widehat{\mathbf{B}}_i'||_* \\
&= (\lambda_B^{(k)} - \frac{1}{c_y} \sum_{i \in \mathcal{I}_k} \lambda_B^{(i)})||\widehat{\mathbf{B}}_k||_* \\
&\geq 0
\end{aligned}
$$

Now assume a violation of condition 2, wherein $\lambda_B^{(k)} \geq \frac{1}{c_{sy}} \sum_{i \in \mathcal{I}_k} \lambda_S^{(i)}||\mathbf{Y}_{\cdot}^{(k)}||_*$. Let $\widehat{\mathbf{B}}_k \mathbf{Y}_{\cdot}^{(k)} = \sum_{i \in \mathcal{I}_k} \widehat{\mathbf{S}}_{\cdot}^{(i)'}$, where $\widehat{\mathbf{S}}_{\cdot}^{(i)'}$ contains the blocks of $\widehat{\mathbf{B}}_k \mathbf{Y}_{\cdot}^{(k)}$ corresponding to $\mathbf{C}_Y[\cdot, i]$ and $\mathbf{0}$ otherwise. The choice of $\widehat{\mathbf{S}}_{\cdot}^{(i)'}$ is unique that $\widehat{\mathbf{S}}_{\cdot}^{(i)'} = \frac{1}{c_{sy}}\widehat{\mathbf{B}}_k \mathbf{Y}_{S\cdot}^{(i)}$, where $\mathbf{Y}_{S\cdot}^{(i)} = [\mathbf{Y}_{S1}^{(i)}, \mathbf{Y}_{S2}^{(i)}, ..., \mathbf{Y}_{S\mathbf{J}}^{(i)}]$ with $\mathbf{Y}_{Sj}^{(i)} = \begin{cases} \mathbf{0}_{q \times n_j} & \text{if } \mathbf{C}_S[j, i] = 0 \\ \mathbf{Y}_j & \text{if } \mathbf{C}_S[j, i] = 1 \end{cases}$ for all $i \in \mathcal{I}_k$. Since $\mathbf{Y}_{S\cdot}^{(i)}$ is gained by setting some blocks of $\mathbf{Y}_{\cdot}^{(k)}$ to be zero, $||\widehat{\mathbf{S}}_{\cdot}^{(i)'}||_* = ||\frac{1}{c_{sy}}\widehat{\mathbf{B}}_k \mathbf{Y}_{S\cdot}^{(i)}||_* \leq ||\frac{1}{c_{sy}}\widehat{\mathbf{B}}_k \mathbf{Y}_{\cdot}^{(k)}||_*$. Consider a minimizer $\{\widetilde{\mathbf{B}}_k\}_{k=1}^{K}, \{\widetilde{\mathbf{S}}_{\cdot}^{(l)}\}_{l=1}^{L}$, where $\widetilde{\mathbf{B}}_k = \mathbf{0}$, $\widetilde{\mathbf{S}}_i = \widehat{\mathbf{S}}_{\cdot}^{(i)} + \widehat{\mathbf{S}}_{\cdot}^{(i)'}$, $\forall i \in \mathcal{I}_k$, and all other $\mathbf{B}, \mathbf{S}$ estimates are equal. Then, by the triangle inequality,

$$
\begin{aligned}
f(\{\widehat{\mathbf{B}}_k\}_{k=1}^{K}, \{\widehat{\mathbf{S}}_{\cdot}^{(l)}\}_{l=1}^{L}) - f(\{\widetilde{\mathbf{B}}_k\}_{k=1}^{K}, \{\widetilde{\mathbf{S}}_{\cdot}^{(l)}\}_{l=1}^{L}) &= \lambda_B^{(k)}||\widehat{\mathbf{B}}_k||_* + \sum_{i \in \mathcal{I}_k} \lambda_S^{(i)}||\widehat{\mathbf{S}}_{\cdot}^{(i)}||_* - \sum_{i \in \mathcal{I}_k} \lambda_S^{(i)}||\widehat{\mathbf{S}}_{\cdot}^{(i)} + \widehat{\mathbf{S}}_{\cdot}^{(i)'}||_* \\
&\geq \lambda_B^{(k)}||\widehat{\mathbf{B}}_k||_* + \sum_{i \in \mathcal{I}_k} \lambda_S^{(i)}||\widehat{\mathbf{S}}_{\cdot}^{(i)}||_* - \sum_{i \in \mathcal{I}_k} \lambda_S^{(i)}(||\widehat{\mathbf{S}}_{\cdot}^{(i)}||_* + ||\widehat{\mathbf{S}}_{\cdot}^{(i)'}||_*) \\
&= \lambda_B^{(k)}||\widehat{\mathbf{B}}_k||_* - \sum_{i \in \mathcal{I}_k} \lambda_S^{(i)}||\widehat{\mathbf{S}}_{\cdot}^{(i)'}||_* \\
&\geq \lambda_B^{(k)}||\widehat{\mathbf{B}}_k||_* - \sum_{i \in \mathcal{I}_k} \lambda_S^{(i)}||\frac{1}{c_{sy}}\widehat{\mathbf{B}}_k \mathbf{Y}_{\cdot}^{(k)}||_* \\
&\geq \lambda_B^{(k)}||\widehat{\mathbf{B}}_k||_* - \sum_{i \in \mathcal{I}_k} \lambda_S^{(i)}||\frac{1}{c_{sy}}\widehat{\mathbf{B}}_k||_*||\mathbf{Y}_{\cdot}^{(k)}||_* \\
&= (\lambda_B^{(k)} - \frac{1}{c_{sy}} \sum_{i \in \mathcal{I}_k} \lambda_S^{(i)}||\mathbf{Y}_{\cdot}^{(k)}||_*)||\widehat{\mathbf{B}}_k||_* \\
&\geq 0
\end{aligned}
$$

The proof for condition 3 and 4 is similar to arguments made in (Lock et al., 2022). ∎

## B.3 Proof of Proposition 2

*Proof.* Since the rows of $\mathbf{Y}$ are linear independent, the projection onto the space spanned by rows of $\mathbf{Y}, \mathcal{R}(\mathbf{Y})$, is $\mathbf{Y}^T(\mathbf{Y}\mathbf{Y}^T)^{-1}\mathbf{Y} = \mathbf{Y}^T\mathbf{Y}$. By decomposing $\mathbf{X}$ onto $\mathcal{R}(\mathbf{Y})$ and $\mathcal{R}(\mathbf{I} - \mathbf{Y})$,

we have

$$||\mathbf{X} - \mathbf{AY}||_F^2 = ||\mathbf{X}(\mathbf{I} - \mathbf{Y}^T\mathbf{Y}) + \mathbf{XY}^T\mathbf{Y} - \mathbf{AY}||_F^2$$
$$= ||\mathbf{X}(\mathbf{I} - \mathbf{Y}^T\mathbf{Y})||_F^2 + ||\mathbf{XY}^T\mathbf{Y} - \mathbf{AY}||_F^2 + tr[2(\mathbf{I} - \mathbf{Y}^T\mathbf{Y})\mathbf{X}^T(\mathbf{XY}^T - \mathbf{A})\mathbf{Y}]$$
$$= ||\mathbf{X}(\mathbf{I} - \mathbf{Y}^T\mathbf{Y})||_F^2 + ||\mathbf{XY}^T\mathbf{Y} - \mathbf{AY}||_F^2 + tr[2\mathbf{Y}(\mathbf{I} - \mathbf{Y}^T\mathbf{Y})\mathbf{X}^T(\mathbf{XY}^T - \mathbf{A})]$$
$$= ||\mathbf{X}(\mathbf{I} - \mathbf{Y}^T\mathbf{Y})||_F^2 + ||\mathbf{XY}^T\mathbf{Y} - \mathbf{AY}||_F^2$$

The second equation comes from the cyclic property of trace. Construct an orthogonal matrix $\mathbf{Q} = [\mathbf{Y}, \mathbf{Y}^*]$ where the rows of $\mathbf{Y}^*$ give an orthonormal basis for $\mathcal{R}(\mathbf{I} - \mathbf{Y})$, i.e. $\mathbf{YY}^{*T} = \mathbf{0}$. Therefore,

$$||\mathbf{XY}^T\mathbf{Y} - \mathbf{AY}||_F^2 = tr[(\mathbf{XY}^T\mathbf{Y} - \mathbf{AY})^T(\mathbf{XY}^T\mathbf{Y} - \mathbf{AY})]$$
$$= tr[(\mathbf{XY}^T\mathbf{Y} - \mathbf{AY})^T(\mathbf{XY}^T\mathbf{Y} - \mathbf{AY})\mathbf{Q}^T\mathbf{Q}]$$
$$= ||\mathbf{XY}^T\mathbf{YQ}^T - \mathbf{AYQ}^T||_F^2$$
$$= ||\mathbf{XY}^T\mathbf{Y}\begin{bmatrix} \mathbf{Y}^T \\ \mathbf{Y}^{*T} \end{bmatrix} - \mathbf{AY}\begin{bmatrix} \mathbf{Y}^T \\ \mathbf{Y}^{*T} \end{bmatrix}||_F^2$$
$$= ||\mathbf{XY}^T\begin{bmatrix} \mathbf{I} \\ \mathbf{0} \end{bmatrix} - \mathbf{A}\begin{bmatrix} \mathbf{I} \\ \mathbf{0} \end{bmatrix}||_F^2$$
$$= ||\mathbf{XY}^T - \mathbf{A}||_F^2$$

Combining these results, we rewrite

$$\min_{\mathbf{A}}\{\frac{1}{2}||\mathbf{X} - \mathbf{AY}||_F^2 + \lambda||\mathbf{A}||_*\} = \min_{\mathbf{A}}\{\frac{1}{2}||\mathbf{X}(\mathbf{I} - \mathbf{Y}^T\mathbf{Y})||_F^2 + \frac{1}{2}||\mathbf{XY}^T - \mathbf{A}||_F^2 + \lambda||\mathbf{A}||_*\}$$
$$= \min_{\mathbf{A}}\{\frac{1}{2}||\mathbf{XY}^T - \mathbf{A}||_F^2 + \lambda||\mathbf{A}||_*\} + \frac{1}{2}||\mathbf{X}(\mathbf{I} - \mathbf{Y}^T\mathbf{Y})||_F^2$$

Apply Lemma 1 and we get the first desired result. The second desired result follows immediately due to $||\mathbf{A}||_* = ||\mathbf{AY}||_*$ for $\mathbf{YY}^T = \mathbf{I}$. Let SVD of $\mathbf{A}$ be $\mathbf{U}_A\mathbf{D}_A\mathbf{V}_A^T$ and $\mathbf{V}_A' = \mathbf{Y}^T\mathbf{V}_A$. Since $\mathbf{V}_A'^T\mathbf{V}_A' = \mathbf{V}_A^T\mathbf{YY}^T\mathbf{V}_A = \mathbf{V}_A^T\mathbf{V}_A = \mathbf{I}$, $\mathbf{U}_A\mathbf{D}_A\mathbf{V}_A'^T$ is the SVD of $\mathbf{AY}$. $\mathbf{A}$ and $\mathbf{AY}$ share the same singular values so that their nuclear norms are the same. ∎

## B.4   Proof of Proposition 4

*Proof.* We start with the special case when $\mathbf{Y}$ is an identity matrix with $q = n$ and $\mathbf{X} = \mathbf{B} + \frac{1}{\sqrt{n}}\mathbf{E}$. It follows directly from (Shabalin and Nobel, 2013) that

$$\sigma_j(\mathbf{X}) \xrightarrow{P} \begin{cases} \sqrt{1 + \sigma_j^2(\mathbf{B}) + c + \dfrac{c}{\sigma_j^2(\mathbf{B})}}, & \text{if } \sigma_j(\mathbf{B}) > \sqrt[4]{c} \\ 1 + \sqrt{c}, & \text{if } \sigma_j(\mathbf{B}) \leq \sqrt[4]{c} \end{cases}$$

. Note $s(\sigma_j(\mathbf{B})) = \sqrt{1 + \sigma_j^2(\mathbf{B}) + c + \frac{c}{\sigma_j^2(\mathbf{B})}}$ is a monotonic increasing function when $\sigma_j(\mathbf{B}) > \sqrt[4]{c}$. Therefore, $min\{s(\sigma_j(\mathbf{B}))\} > s(\sqrt[4]{c}) = 1 + \sqrt{c}$.

Now consider the more general case that $\mathbf{X}_{m \times n} = \mathbf{B}_{m \times q}\mathbf{Y}_{q \times n} + \mathbf{E}_{m \times n}$. Due to semi-orthogonality of $\mathbf{Y}$, $\mathbf{XY}^T = \mathbf{B} + \mathbf{EY}^T$. Since $\mathbf{E}$ is of matrix normal distribution $\mathcal{MN}_{m,n}(\mathbf{0}_{m \times n}, \mathbf{I}_{m \times m}, \mathbf{I}_{n \times n})$ and $\mathbf{Y}_{n \times q}^T$ is a linear transformation of full rank $q \leq n$, we have

$$\mathbf{EY}^T \sim \mathcal{MN}_{m,q}(\mathbf{0Y}^T, \mathbf{I}_{m \times m}, \mathbf{YIY}^T) = \mathcal{MN}_{m,q}(\mathbf{0}_{m \times q}, \mathbf{I}_{m \times m}, \mathbf{I}_{q \times q}).$$

Recognize that entries of $\mathbf{E}\mathbf{Y}^T$ are still independent normal. Denote $\bar{\mathbf{E}} = \mathbf{E}\mathbf{Y}^T$ and $\bar{\mathbf{X}} = \mathbf{X}\mathbf{Y}^T$. The original question becomes $\bar{\mathbf{X}}_{m \times q} = \mathbf{B}_{m \times q} + \bar{\mathbf{E}}_{m \times q}$. Applying the result of the special case,

$$\sigma_j(\mathbf{X}\mathbf{Y}^T) = \sigma_j(\bar{\mathbf{X}}) \xrightarrow{P} \begin{cases} \sqrt{1 + \sigma_j^2(\mathbf{B}) + c + \dfrac{c}{\sigma_j^2(\mathbf{B})}}, & \text{if } \sigma_j(\mathbf{B}) > \sqrt[4]{c} \\ 1 + \sqrt{c}, & \text{if } \sigma_j(\mathbf{B}) \leq \sqrt[4]{c}. \end{cases}$$

Following the aforementioned reasoning, we have the conclusion that $\sigma_j(\mathbf{X}\mathbf{Y}^T)$ will converge to a number larger than $1 + \sqrt{c}$ as $\sigma_j(\mathbf{B}) > \sqrt[4]{c}$. ∎

## B.5 Proposition 5 and its proof

**Proposition 5.** *For any semi-orthogonal matrix $\mathbf{Y}$ such that $\mathbf{Y}^T\mathbf{Y} = \mathbf{I}$, if the optimization problems $\min_{\mathbf{B}}\{\frac{1}{2}||\mathbf{X} - \mathbf{B}\mathbf{Y}||_F^2 + \lambda||\mathbf{B}||_*\}$ and $\min_{\mathbf{S}}\{\frac{1}{2}||\mathbf{X} - \mathbf{S}||_F^2 + \lambda||\mathbf{S}||_*\}$ have their optimal solutions as $\widehat{\mathbf{B}}$ and $\widehat{\mathbf{S}}$ respectively, then $\widehat{\mathbf{S}} = \widehat{\mathbf{B}}\mathbf{Y}$.*

*Proof.* Since orthogonal transformation preserves Frobenius norm, we have $||\mathbf{X} - \mathbf{B}\mathbf{Y}||_F = ||\mathbf{X}\mathbf{Y}^T - \mathbf{B}||_F$. Then,

$$\min_{\mathbf{B}}\{\frac{1}{2}||\mathbf{X} - \mathbf{B}\mathbf{Y}||_F^2 + \lambda||\mathbf{B}||_*\} = \min_{\mathbf{B}}\{\frac{1}{2}||\mathbf{X}\mathbf{Y}^T - \mathbf{B}||_F^2 + \lambda||\mathbf{B}||_*\}.$$

Let SVD of $\mathbf{X}$ to be $\mathbf{U}\mathbf{D}\mathbf{V}^T$. Since $(\mathbf{Y}\mathbf{V})^T(\mathbf{Y}\mathbf{V}) = \mathbf{V}^T\mathbf{Y}^T\mathbf{Y}\mathbf{V} = \mathbf{V}^T\mathbf{V} = \mathbf{I}$, we have $\mathbf{U}\mathbf{D}(\mathbf{Y}\mathbf{V})^T = \mathbf{U}\mathbf{D}\mathbf{V}^T\mathbf{Y}^T = \mathbf{X}\mathbf{Y}^T$ serving as the SVD of $\mathbf{X}\mathbf{Y}^T$. Applying Lemma 1, the solution for $\min_{\mathbf{B}}\{\frac{1}{2}||\mathbf{X}\mathbf{Y}^T - \mathbf{B}||_F^2 + \lambda||\mathbf{B}||_*\}$ is $\widehat{\mathbf{B}} = \mathbf{U}\widetilde{\mathbf{D}}(\mathbf{Y}\mathbf{V})^T$ and the solution for $\min_{\mathbf{S}}\{\frac{1}{2}||\mathbf{X} - \mathbf{S}||_F^2 + \lambda||\mathbf{S}||_*\}$ is $\widehat{\mathbf{S}} = \mathbf{U}\widetilde{\mathbf{D}}\mathbf{V}^T$. Therefore, $\widehat{\mathbf{S}} = \widehat{\mathbf{B}}\mathbf{Y}$.

∎

If $\mathbf{Y}_{\cdot}^{(k)}$ is semi-orthogonal with $\mathbf{Y}_{\cdot}^{(k)T}\mathbf{Y}_{\cdot}^{(k)} = \mathbf{I}_n(q \geq n)$ and contains information from exactly the same cohorts as $\mathbf{S}^{(l)}$ (i.e., $\mathbf{C}_Y[\cdot, k] = \mathbf{C}_S[\cdot, l]$), the estimation process cannot differentiate $\widehat{\mathbf{B}}_k\mathbf{Y}_{\cdot}^{(k)}$ from $\widehat{\mathbf{S}}^{(l)}$. Even without overlapping modules, the newly proposed model reverts to the unsupervised model 3 described in Section 3, which only comprises $\{\widehat{\mathbf{S}}_{\cdot}^{(l)}\}_{l=1}^L$. Thus, the loss with semi-orthogonal regressors $\mathbf{Y} : q \times n$ is only valid if $q < n$. This finding aligns with our real-world problem of interest: in reality, we have only a few somatic mutations, the number of which is less than the number of patients in the study.

# Appendix C  Scaling and orthogonalization

In reality the original $\mathbf{Y}$ is not orthogonal. In practice we scale $\mathbf{X}$ and orthogonalize $\mathbf{Y}$ prior to optimization. We first center each row of $\mathbf{X}_{\cdot}$ to have mean 0. In order to satisfy the standard normal noise requirement, we estimate the error variance for $\mathbf{X}_{\cdot}$ by using the median absolute deviation estimator from (Gavish and Donoho, 2017). The estimated variance of $\mathbf{X}_{\cdot}$ is denoted as $\hat{\sigma}^2$. Then, we use $\mathbf{X}_{\cdot}/\hat{\sigma}$ as the final data matrix for optimization, which has residual variance approximately 1. We further orthogonalize the columns of $\mathbf{Y}$ via SVD prior to optimization, and transform the solution back to the original covariate space afterward. Here, we describe two scaling approaches for $\mathbf{Y}$: one standardizing the original covariates with out orthogonalization, and another rotating the covariate space so that it is orthogonal. In the next section, we show simulation results that illustrate the difference between these two approaches.

Instead of estimating $\mathbf{B}$ based on original $\mathbf{Y}_{\cdot}^{(k)}$, we first center each covariate in $\mathbf{Y}_{\cdot}^{(k)}$ and scale each covariate to have variance 1 and then do the optimization. For each $k = 1, ..., K$, the detailed steps are:

1. Center each covariate in $\mathbf{Y}^{(k)}_{\cdot}$ and calculate the square root of row sums of squared entry of $\mathbf{Y}^{(k)}_{\cdot}$, denoted as $t^{(k)}_1, t^{(k)}_2, ..., t^{(k)}_q$.

2. Construct a diagonal matrix $\mathbf{\Sigma}_{\boldsymbol{nk}} = diag(t^{(k)}_1, ..., t^{(k)}_q)$.

3. Use $\mathbf{Y}^{(k)}_{\cdot scaled} = \mathbf{\Sigma}_{\boldsymbol{nk}}^{-1}\mathbf{Y}^{(k)}_{\cdot}$ in optimization.

4. Denote the estimation with respect to $\mathbf{Y}^{(k)}_{\cdot scaled}$ as $\widehat{\mathbf{B}}_{k,scaled}$. Then, our final estimation for $\mathbf{B}_k$ on its original scale is $\widehat{\mathbf{B}}_k = \widehat{\mathbf{B}}_{k,scaled}\mathbf{\Sigma}_{\boldsymbol{nk}}$.

In order to further reduce collinearity among covariates, after centerization we may instead choose to apply orthogonalization to $\mathbf{Y}^{(k)}_{\cdot}$ rather than standardization. For each $k = 1, ..., K$, the detailed steps are:

1. Obtain the SVD of $\mathbf{Y}^{(k)}_{\cdot}$, i.e. $\mathbf{Y}^{(k)}_{\cdot} = \mathbf{U}_k\mathbf{D}_k\mathbf{V}_k^T$.

2. Use $\mathbf{Y}^{(k)}_{\cdot orth} = \mathbf{V}_k^T$ in optimization.

3. Denote the estimation with respect to $\mathbf{Y}^{(k)}_{\cdot orth}$ as $\widehat{\mathbf{B}}_{k,orth}$. Then, our final estimation for $\mathbf{B}_k$ on its original scale is $\widehat{\mathbf{B}}_k = \widehat{\mathbf{B}}_{k,orth}\mathbf{U}_k\mathbf{D}_k$.

There are other ways to orthogonalize, such as Gram-Schmidt method. No matter standardization or orthogonalization, we record the transforming matrix from the original $\mathbf{Y}^{(k)}_{\cdot}$ to new $\mathbf{Y}^{(k)}_{\cdot scaled/orth}$ for correction of estimated $\mathbf{B}_k$ at the final stage.

# Appendix D  More details on the simulations

## D.1  Complete data generation

Here we describe the complete process to generate data in Section 6.2.

1. Every entry in each $\mathbf{Y}_j, j = 1, ..., J = 30$ is drawn independently from a standard normal distribution. By default, the variance of each feature is 1. Construct one global shared $\mathbf{Y}^{(1)}_{\cdot} = [\mathbf{Y}_1, ..., \mathbf{Y}_{30}]$ and 30 individual $\mathbf{Y}^{(k)}_{\cdot} = [\mathbf{0}, ..., \mathbf{Y}_{k-1}, ..., \mathbf{0}], k = 2, ..., 31$ (only the $k-1$th submatrix is non-zero).

2. For the $K = 31$ modules of $\mathbf{Y}^{(k)}_{\cdot}$, generate $\mathbf{B}_k = \mathbf{U}^{(k)}_B\mathbf{V}^{(k)T}_B/sd(\mathbf{U}^{(k)}_B\mathbf{V}^{(k)T}_B\mathbf{Y}^{(k)}_{\cdot}) * \sqrt{n_k/n}$, where $n_k$ is the number of samples in module $k$. Each entry of $\mathbf{U}^{(k)}_B : p \times r, \mathbf{V}^{(k)}_B : q \times r$ comes from standard Normal distribution.

3. For a number of $L = 31$ modules of $\mathbf{S}^{(l)}_{\cdot}$ involved, draw one global score matrix $\mathbf{V}^{(1)}_S : n \times r$ and 30 individual score matrices $\mathbf{V}^{(l)}_S = [\mathbf{0}, ..., \mathbf{V}_{l-1}, ..., \mathbf{0}], l = 2, ..., 31$, where each entry of $\mathbf{V}_1, ..., \mathbf{V}_{30}$ comes from standard Normal. Generate $\mathbf{S}^{(l)}_{\cdot} = \mathbf{U}^{(l)}_S\mathbf{V}^{(l)T}_S/sd(\mathbf{U}^{(l)}_S\mathbf{V}^{(l)T}_S) * \sqrt{n_l/n}$, where each entry of $\mathbf{U}^{(l)}_S : p \times r$ comes from standard Normal distribution and $n_l$ is the number of samples in module $l$.

4. Draw each entry of $\mathbf{E}_{\cdot}$ from a standard normal distribution.

5. Generate $\mathbf{X}_{\cdot} = a * \mathbf{B}_1\mathbf{Y}^{(1)}_{\cdot} + b * \mathbf{S}^{(1)}_{\cdot} + c * \sum_{k=2}^{31}\mathbf{B}_k\mathbf{Y}^{(k)}_{\cdot} + d * \sum_{l=2}^{31}\mathbf{S}^{(l)}_{\cdot} + \mathbf{E}_{\cdot\cdot}$. The letters $a, b, c, d$ are constant for signal size. For instance, scenario $(a)$, $a = \sqrt{10}$ and the remaining equal 1.

## D.2   Computation time

Table 6 demonstrates the computation time for all methods in missing data imputation for the TCGA data. The computation time of non-missing optimization and missing data imputation are similar. The proposed maRRR method consumes the most time since it considers both covariate effects and auxiliary modules. Generally, the computation time is proportional to the number of modules. But it will vary because of different number of cohorts within one module. Algorithm 1 requests more computation time when the true rank in some module is large. In general, both algorithms output similar RSE so that the one takes less computation time is used in the simulation. The case of missing columns takes slightly less computation time than the other two cases. This results from that it is uninformative when missing certain subjects. It may lead to all-zero imputation and this explains why aRRR for 30 separate modules takes significantly less time than the other cases.

Table 6: Computation time (in seconds) for our proposed methods and all other comparison methods in missing data imputation (based on algorithm 1 and 2 respectively) for the TCGA real data application. Each other method is based on 30 epochs. NA represents that missing data imputation based on Algorithm 2 does not work for "NNreg" and "NNapprox". "50 modules" means that the method is based on 50 detected modules. "1 all-shared + 30 separate" means that the method is based on one all-shared module and 30 separate modules (in total 31 modules). For method name representations, please refer to Section 6.2 in the main context.

| Method | Algorithm 1 | | | | Algorithm 2 | | | |
| | missing entries | missing columns | missing rows | average | missing entries | missing columns | missing rows | average |
|---|---|---|---|---|---|---|---|---|
| maRRR 50 modules | 6251.0 | 5345.2 | 6442.4 | 6012.9 | 3725.3 | 3264.7 | 3754.8 | 3581.6 |
| BIDIFAC+, 50 modules | 5000.7 | 4283.4 | 5198.2 | 4827.4 | 1140.7 | 1058.2 | 1161.3 | 1120.1 |
| mRRR, 50 modules | 1566.7 | 1327.9 | 1576.9 | 1490.5 | 2905.9 | 2519.9 | 2873.3 | 2766.4 |
| maRRR, 1 all-shared + 30 separate | 1925.6 | 1623.3 | 1962.9 | 1837.3 | 2276.5 | 2085.1 | 2303.6 | 2221.7 |
| BIDIFAC+, 1 all-shared + 30 separate | 1109.9 | 967.1 | 1124.9 | 1067.3 | 705.2 | 684.8 | 714.9 | 701.6 |
| mRRR, 1 all-shared + 30 separate | 938.2 | 823.7 | 939.5 | 900.4 | 1711.2 | 1556.3 | 1729.0 | 1665.5 |
| aRRR, one all-shared | 617.9 | 524.8 | 637.8 | 593.5 | 75.1 | 67.6 | 74.7 | 72.5 |
| aRRR, 30 separate | 1246.3 | 68.5 | 1249.0 | 854.6 | 2196.8 | 123.4 | 2195.6 | 1505.3 |
| NNreg, one all-shared | 2.2 | 2.2 | 5.5 | 3.3 | NA | NA | NA | NA |
| NNapprox, one all-shared | 116.7 | 49.6 | 585.0 | 250.5 | NA | NA | NA | NA |
| NNreg, 30 separate | 4.5 | 18.5 | 21.5 | 14.9 | NA | NA | NA | NA |
| NNapprox, 30 separate | 60.4 | 22.6 | 256.2 | 113.1 | NA | NA | NA | NA |

## D.3 Additional simulation to assess aRRR and maRRR

Here we discuss an additional simulation study to assess aspects of aRRR and maRRR, including the effects of orthogonalizing Y and the recovery of the underlying ranks of the true structure. The simulation covers scenarios in which the original explanatory data matrices $\mathbf{Y}_j, j = 1, ..., J$ are orthogonal or not. The process, used to generate the complete data for this section, is as follows:

1. Every entry in each $\mathbf{Y}_j, j = 1, ..., J$ is drawn independently from a standard normal distribution. By default, the variance of each feature is 1. Denote the standard deviation for each $\mathbf{Y}_j, j = 1, ..., J$ as $sd(\mathbf{Y}_j), j = 1, ..., J$.

2. If we aim to orthogonalize some $\mathbf{Y}_j, j = 1, ..., J$, denote the transpose of the right orthogonal matrix from singular value decomposition of $\mathbf{Y}_j$ as $\mathbf{V}_{Y_j}$. Calculate the orthogonal version of $\mathbf{Y}_j, j = 1, ..., J$ by $\mathbf{V}_{Y_j}/sd(\mathbf{Y}_j), j = 1, ..., J$ in order to keep the original scale of variation.

3. For a number of $K$ modules of covariate effects involved, construct $\mathbf{Y}_.^{(k)}, k = 1, ..., K$. Generate $\mathbf{B}_k = \mathbf{U}_B^{(k)}\mathbf{V}_B^{(k)T}$, where each entry of $\mathbf{U}_B^{(k)} : q \times r, \mathbf{V}_B^{(k)} : n_k \times r$ comes from standard Normal distribution.

4. For a number of $L$ modules of $\mathbf{S}^{(l)}$ involved, generate $\mathbf{S}_.^{(l)} = \mathbf{U}_S^{(l)}\mathbf{V}_S^{(l)T}$, where each entry of $\mathbf{U}_S^{(l)} : p \times r$ and each non-zero entry of $\mathbf{V}_S^{(l)} : q \times r$ comes from standard Normal distribution.

5. Draw each entry of $\mathbf{E}_.$ from a standard normal distribution.

6. Generate $\mathbf{X}_. = \sum_{k=1}^{K} \mathbf{B}_k \mathbf{Y}_.^{(k)} + \sum_{l=1}^{L} \mathbf{S}_.^{(l)} + \mathbf{E}_{..}$

Besides mean squared error (mse) as a metric to assess accuracy, we want to understand how low-rank structures for $\mathbf{B}, \mathbf{S}$ are uncovered is estimated under different orthogonality. Therefore, we define the "rank sum ratio" as follows:

$$\frac{\sum_{i=1}^{r_B} \lambda_i(\widehat{\mathbf{B}})}{\sum_{i=r_B+1}^{r_{B,upper}} \lambda_i(\widehat{\mathbf{B}})}, \frac{\sum_{i=1}^{r_S} \lambda_i(\widehat{\mathbf{S}})}{\sum_{i=r_S+1}^{r_{S,upper}} \lambda_i(\widehat{\mathbf{S}})}$$

where $r_B$ is the true rank of $\mathbf{B}$ and $r_{B,upper}$ is the upper bound of $\mathbf{B}$ specified in estimation, both similarly defined for $\mathbf{S}$. The smaller the rank sum ratio is, the lower the rank structure achieves.

Table 7 and 8 shows the MSE and rank sum ratio for aRRR and maRRR simulations respectively. We consider one cohort for aRRR and two cohorts for maRRR. In general, orthogonaliztion before optimization and in data generation achieve similar results. Compared with only standardization of $\mathbf{Y}$, orthogonaliztion of $\mathbf{Y}$ leads to less MSE and more accurate rank estimations. After orthogonaliztion, the proposed methods still overestimate the true rank of covariate-related signal, but it is not severe as the rank sum ratio is below 0.01. Therefore, orthogoalization of $\mathbf{Y}$ is helpful for optimization.

Table 7: Mean squared error(MSE) for aRRR under different scenarios of orthogonality, true rank and signal size. Row names: "r_y" refers to true rank of $\mathbf{Y}$ in the data generation process. "sd_YB" and "sd_S" refer to the standard deviation for matrix $\mathbf{BY}$ and $\mathbf{S}$ in generation respectively. "epochs" refers to the number of iterations to converge. "ratio_B" refers to the rank sum ratio defined before. "est_rank_B" counts the number of singular values that large than 0.1. Column names: "no orth" means that the original $\mathbf{Y}$ is not orthogonal and we only standardize it before optimization. "orth_opt" means that we only orthogonalize original "$\mathbf{Y}$" before optimization. "orth_gen" means that the original $\mathbf{Y}$ is orthogonal.

|          | r_y | sd_YB | sd_S | epochs | mse_B | mse_S | est_rank_B | ratio_B |
|----------|-----|-------|------|--------|-------|-------|------------|---------|
| no orth  | 1   | 5     | 0.5  | 65.82  | 0.001 | 0.607 | 1.3        | 0.002   |
| no orth  | 1   | 1     | 1    | 31.27  | 0.038 | 0.225 | 1.39       | 0.012   |
| no orth  | 1   | 0.5   | 5    | 44.82  | 0.172 | 0.012 | 1.48       | 0.037   |
| no orth  | 5   | 5     | 0.5  | 52.36  | 0.007 | 0.629 | 5.04       | 0       |
| no orth  | 5   | 1     | 1    | 34.38  | 0.141 | 0.238 | 4.95       | 0.001   |
| no orth  | 5   | 0.5   | 5    | 41.5   | 0.404 | 0.013 | 4.42       | 0       |
| orth_opt | 1   | 5     | 0.5  | 63.66  | 0.001 | 0.606 | 1.06       | 0       |
| orth_opt | 1   | 1     | 1    | 33.69  | 0.033 | 0.224 | 1.12       | 0.004   |
| orth_opt | 1   | 0.5   | 5    | 47.4   | 0.157 | 0.012 | 1.16       | 0.011   |
| orth_opt | 5   | 5     | 0.5  | 45.41  | 0.006 | 0.626 | 5          | 0       |
| orth_opt | 5   | 1     | 1    | 27.62  | 0.125 | 0.238 | 4.94       | 0       |
| orth_opt | 5   | 0.5   | 5    | 39.7   | 0.372 | 0.013 | 4.42       | 0       |
| orth_gen | 1   | 5     | 0.5  | 63.31  | 0.001 | 0.606 | 1.08       | 0.001   |
| orth_gen | 1   | 1     | 1    | 33.47  | 0.034 | 0.224 | 1.1        | 0.004   |
| orth_gen | 1   | 0.5   | 5    | 47.82  | 0.161 | 0.012 | 1.14       | 0.012   |
| orth_gen | 5   | 5     | 0.5  | 45.42  | 0.006 | 0.624 | 5          | 0       |
| orth_gen | 5   | 1     | 1    | 27.72  | 0.122 | 0.237 | 4.96       | 0       |
| orth_gen | 5   | 0.5   | 5    | 41.2   | 0.369 | 0.013 | 4.34       | 0.001   |

Table 8: Mean squared error(MSE) for maRRR under different scenarios of orthogonality, true rank and signal size. Row names: "r_y" refers to true rank of **Y** in the data generation process. "sd_YB" and "sd_S" refer to the standard deviation for matrix **BY** and **S** in generation respectively. "epochs" refers to the number of iterations to converge. "ratio_B" refers to the rank sum ratio defined before. "est_rank_B" counts the number of singular values that large than 0.1. Column names: "no orth" means that the original **Y** is not orthogonal and we only standardize it before optimization. "orth_opt" means that we only orthogonalize original "**Y**" before optimization. "orth_gen" means that the original **Y** is orthogonal. All results are the average of all **B** or **S**.

| | r_y | sd_YB | sd_S | epochs | mse_B | mse_S | est_rank_B | ratio_B |
|---|---|---|---|---|---|---|---|---|
| no orth | 1 | 2 | 0.2 | 84.56 | 0.008 | 0.953 | 1.493 | 0.008 |
| no orth | 1 | 1 | 1 | 77.72 | 0.037 | 0.175 | 1.873 | 0.024 |
| no orth | 1 | 0.2 | 2 | 88.6 | 0.4 | 0.065 | 1.76 | 0.141 |
| no orth | 5 | 2 | 0.2 | 123.45 | 0.115 | 0.962 | 6.617 | 0.033 |
| no orth | 5 | 1 | 1 | 87.78 | 0.202 | 0.195 | 6.06 | 0.023 |
| no orth | 5 | 0.2 | 2 | 91.41 | 0.735 | 0.065 | 3.273 | 0 |
| orth_opt | 1 | 2 | 0.2 | 82.34 | 0.007 | 0.952 | 1.393 | 0.007 |
| orth_opt | 1 | 1 | 1 | 76.96 | 0.034 | 0.174 | 1.66 | 0.019 |
| orth_opt | 1 | 0.2 | 2 | 90.04 | 0.366 | 0.065 | 1.543 | 0.089 |
| orth_opt | 5 | 2 | 0.2 | 119.61 | 0.114 | 0.955 | 6.727 | 0.036 |
| orth_opt | 5 | 1 | 1 | 86.66 | 0.194 | 0.195 | 6.18 | 0.026 |
| orth_opt | 5 | 0.2 | 2 | 93.63 | 0.702 | 0.065 | 3.313 | 0 |
| orth_gen | 1 | 2 | 0.2 | 82.64 | 0.007 | 0.952 | 1.333 | 0.006 |
| orth_gen | 1 | 1 | 1 | 78.4 | 0.034 | 0.175 | 1.64 | 0.018 |
| orth_gen | 1 | 0.2 | 2 | 89.62 | 0.367 | 0.065 | 1.52 | 0.09 |
| orth_gen | 5 | 2 | 0.2 | 116.33 | 0.109 | 0.955 | 6.743 | 0.035 |
| orth_gen | 5 | 1 | 1 | 86.79 | 0.186 | 0.195 | 6.2 | 0.027 |
| orth_gen | 5 | 0.2 | 2 | 93.07 | 0.692 | 0.065 | 3.437 | 0 |

# Appendix E    More details on the real data analysis

## E.1    Cancer type and mutation details

We provide Table 9 as the summary of all the 30 cancer types and Table 10 as the summary of all the 50 somatic mutations considered in our real data analysis (Section 7).

Table 9: The cancer study abbreviations, sample sizes and study names, sourcing from National Cancer Institute https://gdc.cancer.gov/resources-tcga-users/tcga-code-tables/tcga-study-abbreviations.

| Abbreviation | Sample Size | Study Name |
| --- | --- | --- |
| ACC | 77 | Adrenocortical carcinoma |
| BLCA | 129 | Bladder Urothelial Carcinoma |
| BRCA | 976 | Breast invasive carcinoma |
| CESC | 193 | Cervical squamous cell carcinoma and endocervical adenocarcinoma |
| COAD | 147 | Colon adenocarcinoma |
| ESCA | 184 | Esophageal carcinoma |
| GBM | 150 | Glioblastoma multiforme |
| HNSC | 279 | Head and Neck squamous cell carcinoma |
| KICH | 66 | Kidney Chromophobe |
| KIRC | 415 | Kidney renal clear cell carcinoma |
| KIRP | 161 | Kidney renal papillary cell carcinoma |
| LAML | 170 | Acute Myeloid Leukemia |
| LGG | 283 | Brain Lower Grade Glioma |
| LIHC | 195 | Liver hepatocellular carcinoma |
| LUAD | 230 | Lung adenocarcinoma |
| LUSC | 178 | Lung squamous cell carcinoma |
| OV | 115 | Ovarian serous cystadenocarcinoma |
| PAAD | 150 | Pancreatic adenocarcinoma |
| PCPG | 179 | Pheochromocytoma and Paraganglioma |
| PRAD | 331 | Prostate adenocarcinoma |
| READ | 64 | Rectum adenocarcinoma |
| SARC | 245 | Sarcoma |
| SKCM | 342 | Skin Cutaneous Melanoma |
| STAD | 275 | Stomach adenocarcinoma |
| TGCT | 149 | Testicular Germ Cell Tumors |
| THCA | 400 | Thyroid carcinoma |
| THYM | 119 | Thymoma |
| UCEC | 242 | Uterine Corpus Endometrial Carcinoma |
| UCS | 57 | Uterine Carcinosarcoma |
| UVM | 80 | Uveal Melanoma |

Table 10: Gene labels for 50 somatic mutation data in order.

| Index | Mutation | Index | Mutation | Index | Mutation |
|-------|----------|-------|----------|-------|----------|
| 1 | TP53 | 18 | FAT4 | 35 | PKHD1L1 |
| 2 | TTN | 19 | HMCN1 | 36 | RYR1 |
| 3 | MUC16 | 20 | CSMD1 | 37 | RYR3 |
| 4 | PIK3CA | 21 | MUC5B | 38 | NEB |
| 5 | CSMD3 | 22 | ZFHX4 | 39 | PCDH15 |
| 6 | LRP1B | 23 | FAT3 | 40 | DST |
| 7 | KRAS | 24 | SPTA1 | 41 | MLL3 |
| 8 | RYR2 | 25 | GPR98 | 42 | MLL2 |
| 9 | MUC4 | 26 | PTEN | 43 | MACF1 |
| 10 | FLG | 27 | FRG1B | 44 | DNAH9 |
| 11 | SYNE1 | 28 | AHNAK2 | 45 | BRAF |
| 12 | USH2A | 29 | APOB | 46 | DNAH11 |
| 13 | PCLO | 30 | ARID1A | 47 | DNAH8 |
| 14 | APC | 31 | LRP2 | 48 | CSMD2 |
| 15 | DNAH5 | 32 | XIRP2 | 49 | MUC2 |
| 16 | OBSCN | 33 | Unknown | 50 | ABCA13 |
| 17 | MUC17 | 34 | DMD | | |

## E.2   Data processing and distributions

The pan-cancer RNASeq data, described in Hoadley et al. (2018), were downloaded as the file 'EBPlusPlusAdjustPANCAN_IlluminaHiSeq_RNASeqV2.geneExp.tsv' from https://gdc.cancer.gov/about-data/publications/PanCan-CellOfOrigin [accessed June 23, 2021]. This dataset had undergone preprocessing steps described in Hoadley et al. (2018), including batch correction using an empirical Bayes approach and upper-quartile normalization. These data were further log-transformed (via a $\log(1 + x)$ transformation), and filtered to the 1000 genes with the highest standard deviation after log-transformation. The log-transformed and filtered data were then gene-centered by subtracting the overall mean (across all cancer types) for each gene. The distribution of processed expression values for each cancer type are shown in Figure 3; the distributions are roughly similar and approximately bell-shaped across the different cancer types.

The somatic mutation data were binary, prior to the the scaling described in Section C, where '1' implies there is a somatic mutation in the gene for the given sample and '0' implies there is no somatic mutation. Figure 4 gives the proportion of samples that have a mutation across the 50 genes considered, for each cancer type. Some cancer types have several genes that are frequently mutated, while others have a sparser profile with no frequently mutated genes among those considered.
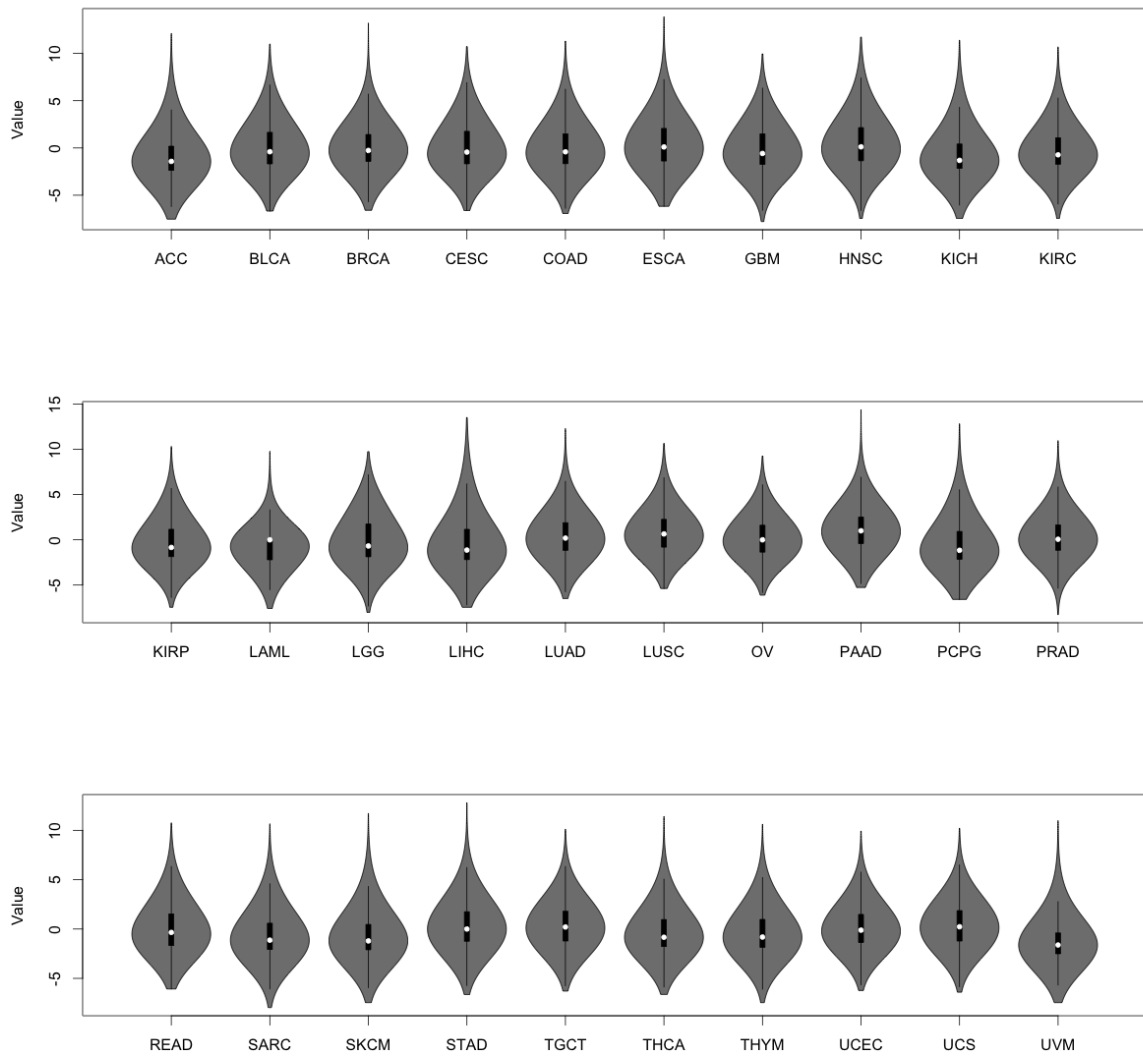
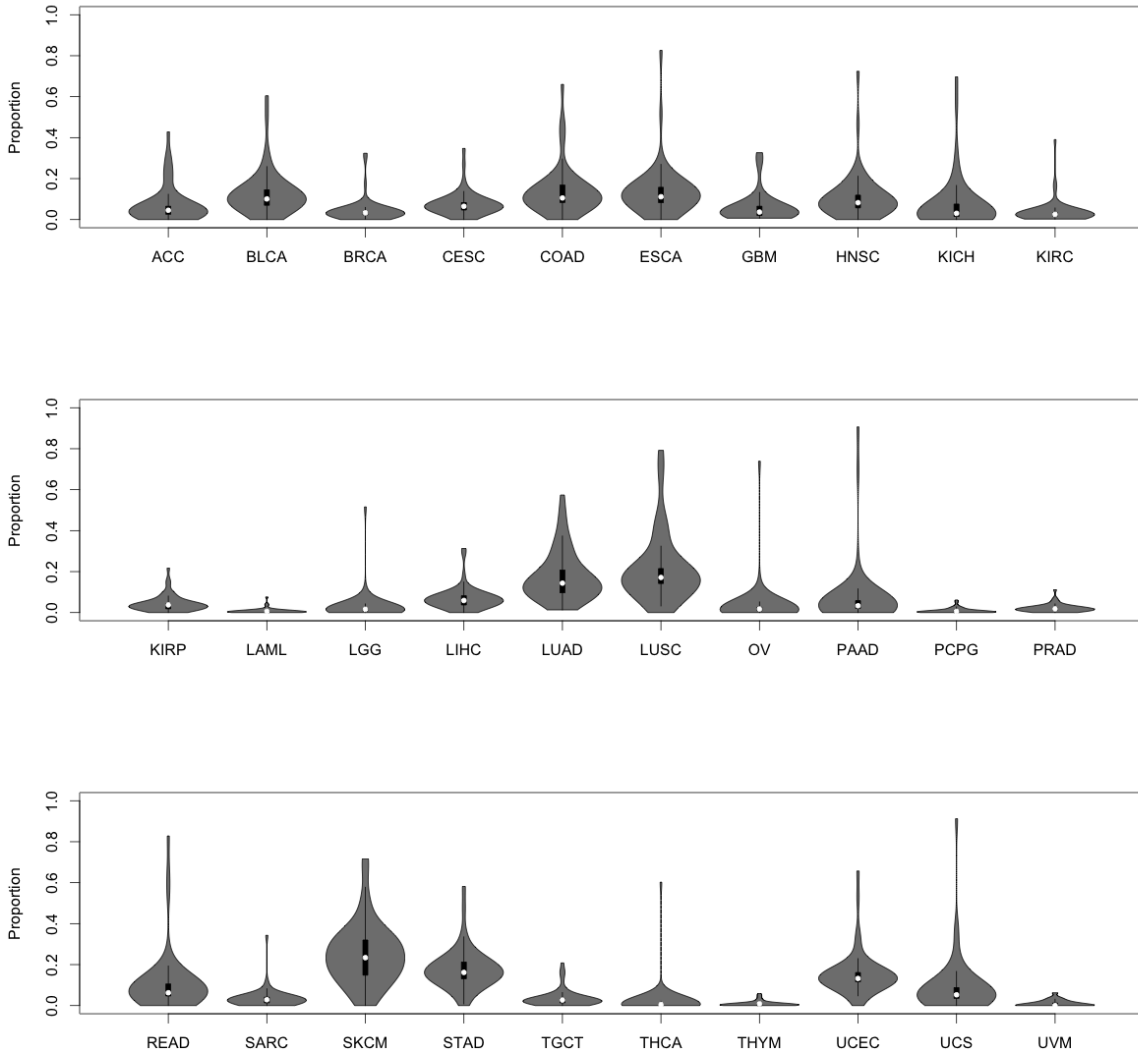Figure 3: Violin plot of normalized expression values for each cancer type.

Figure 4: Violin plot of proportion of samples with a somatic mutation across the 50 genes, for each cancer type.

## E.3 Selection of model parameters

As mentioned in Section 7.2 of the main article, we first apply the optimization with dynamic modules for BIDIFAC+ (Lock et al., 2022), to uncover 50 low-rank modules in $\mathbf{X}_{..}$ BIDIFAC+, noted in Section 3, is the unsupervised version of our proposed model maRRR. The statistical model of BIDIFAC+ is

$$\mathbf{X}_{.} = \sum_{l=1}^{L} \mathbf{S}_{.}^{(l)} + \mathbf{E}_{.}$$

where $\mathbf{S}_{\cdot}^{(l)} = [\mathbf{S}_1^{(l)}, \mathbf{S}_2^{(l)}, ..., \mathbf{S}_J^{(l)}]$, $\mathbf{E}_{\cdot} = [\mathbf{E}_1, \mathbf{E}_2, ..., \mathbf{E}_J]$ and

$$\mathbf{S}_j^{(l)} = \begin{cases} \mathbf{0}_{p \times n_j} & \text{if } \mathbf{C}_S[j,l] = 0 \\ \mathbf{U}_S^{(l)} \mathbf{V}_{Sj}^{(l)T} & \text{if } \mathbf{C}_S[j,l] = 1. \end{cases}$$

The loss objective is

$$\min_{\{\mathbf{S}_{\cdot}^{(l)}\}_{l=1}^{L}} \{\frac{1}{2}||\mathbf{X}_{\cdot} - \sum_{l=1}^{L} \mathbf{S}_{\cdot}^{(l)}||_F^2 + \sum_{l=1}^{L} \lambda_S^{(l)} ||\mathbf{S}_{\cdot}^{(l)}||_*\}.$$

The model aggregates all signals as $\sum_{l=1}^{L} \mathbf{S}_{\cdot}^{(l)}$, yet it doesn't differentiate whether these signals are related to covariates or not. The forward selection process to determine $\mathbf{C}_S$ of BIDIFAC+ initiates with $\mathbf{C}_S[:,l] = \mathbf{0}$ for all $l = 1, ..., L$, and progressively includes cohorts $j$ ($\mathbf{C}_S[j,l] = 1$) to minimize the objective function; see Section 6.3 in (Lock et al., 2022) for complete details. This update occurs iteratively through gradient descent, similar to our Algorithm 2 in Section 5.1 but with dynamic module memberships. This iterative process continues until convergence of loss, at which point the current $\mathbf{C}_S$ is considered the final set of modules. After thorough analysis, we opt for $L = 50$ since including more modules does not significantly enhance our ability to explain variance. A few modules, selected at various $L$ values, contribute substantially to the final model's variance. Consequently, we select the corresponding module indicator matrix $\mathbf{C}_S$ for $L = 50$, thus setting $\mathbf{C}_Y = \mathbf{C}_S$ to effectively partition the variance linked to covariate effects. A comprehensive visualization of final estimated module information is presented in the subsequent heatmap in Figure 5.
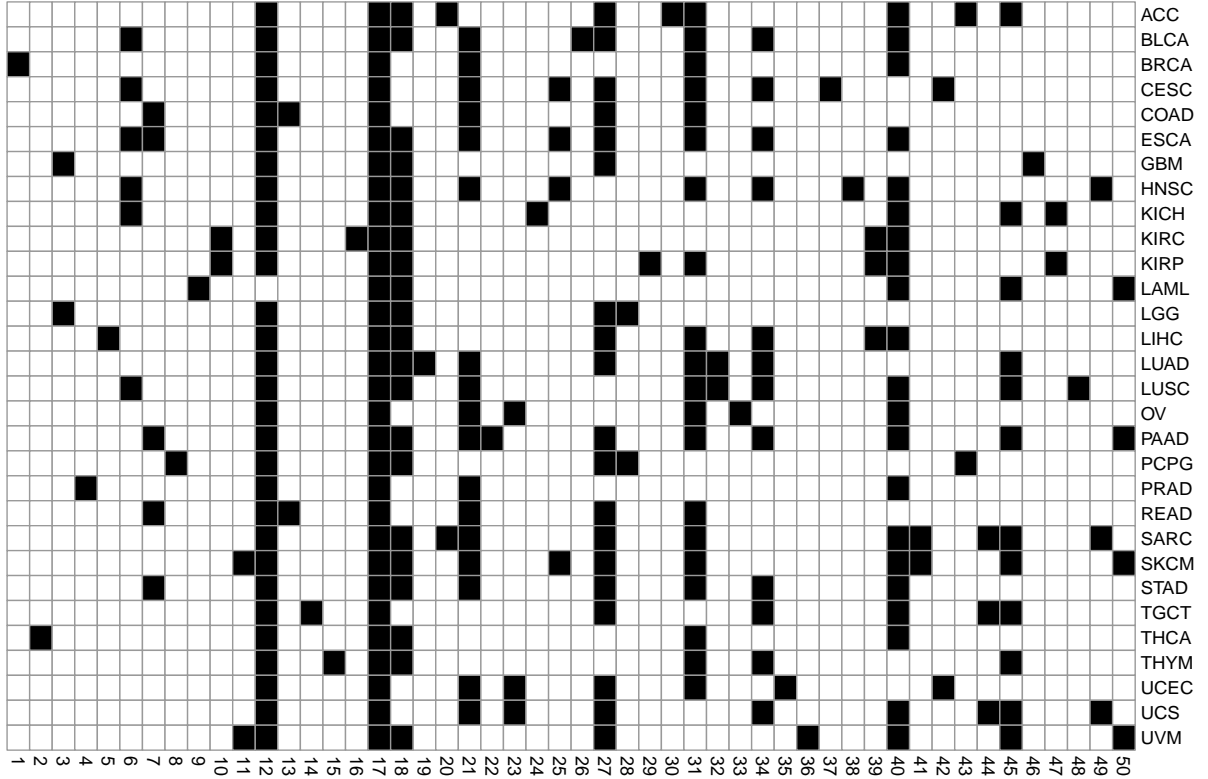


Figure 5: Heatmap for 50 modules detected by BIDIFAC+ model (the chosen $\mathbf{C}_S$) in the TCGA real data application. The row names represent the 30 cohorts while the column names represent the 50 modules. Black grids mean existence of cohorts in the current module.

As outlined in Appendix C, we apply scaling to the standardized $\mathbf{X}$. using the median absolute deviation estimator ((Gavish and Donoho, 2017)) to ensure that the residual variance is approximately 1, i.e., $Var(\mathbf{E}.) = 1$. Drawing from random matrix theorems with a variance of 1 in Section 4, we assign $\lambda_B^{(i)} = \sqrt{1000} + \sqrt{50}, i = 1, ..., 50$ since all $\mathbf{B}_i$ share the same matrix size. Additionally, we set $\lambda_S^{(1)} = \sqrt{1000} + \sqrt{976}$, $\lambda_S^{(2)} = \sqrt{1000} + \sqrt{400}$, ..., and $\lambda_S^{(50)} = \sqrt{1000} + \sqrt{742}$, with the second term representing the square root of the number of samples in the respective module.

In Algorithm 1, we establish a general upper bound for the estimated rank of each $\mathbf{B}_i$ and $\mathbf{S}^{(i)}$ at 20, denoted as $r_{B,upper} = 20$ and $r_{S,upper} = 20$ respectively. Through experimentation, we determined that 20 is the minimum value for Algorithm 1 to match the performance of Algorithm 2. This rank upper bound of 20 is deemed reasonable for low-rank approximations. While the true ranks of most estimates hover around 10, opting for smaller maximum rank values dismisses significant signals. Conversely, larger values don't provide substantial additional information for improving predictions, but rather increase computational complexity.

## E.4 Additional pan-cancer analysis

With the identification of 50 modules and a maximum rank set at 20, maRRR attains an impressive Relative Squared Error (RSE) of 0.184, signifying the substantial capture of variation. Following Section 7.2, the subsequent scatterplot Figure 6 depicts the extent to which mutations account for variance within each module. While the impact of mutations is not overwhelmingly dominant, it is nevertheless significant, aligning with our initial expectations. This outcome is consistent with the extensive genetic information present in gene expressions that is unrelated to mutations. When considered alongside the heatmap of $\mathbf{C}_S$ (Figure 5 in Appendix E.3), it becomes evident that modules such as numbers 2, 6, 7, 13, 22, 23, 41, and 50 exhibit substantial influence from mutations. Remarkably, nearly all of the 30 cancer types make varying appearances within those modules, with the exceptions of ACC, KIRC, and PRAD.
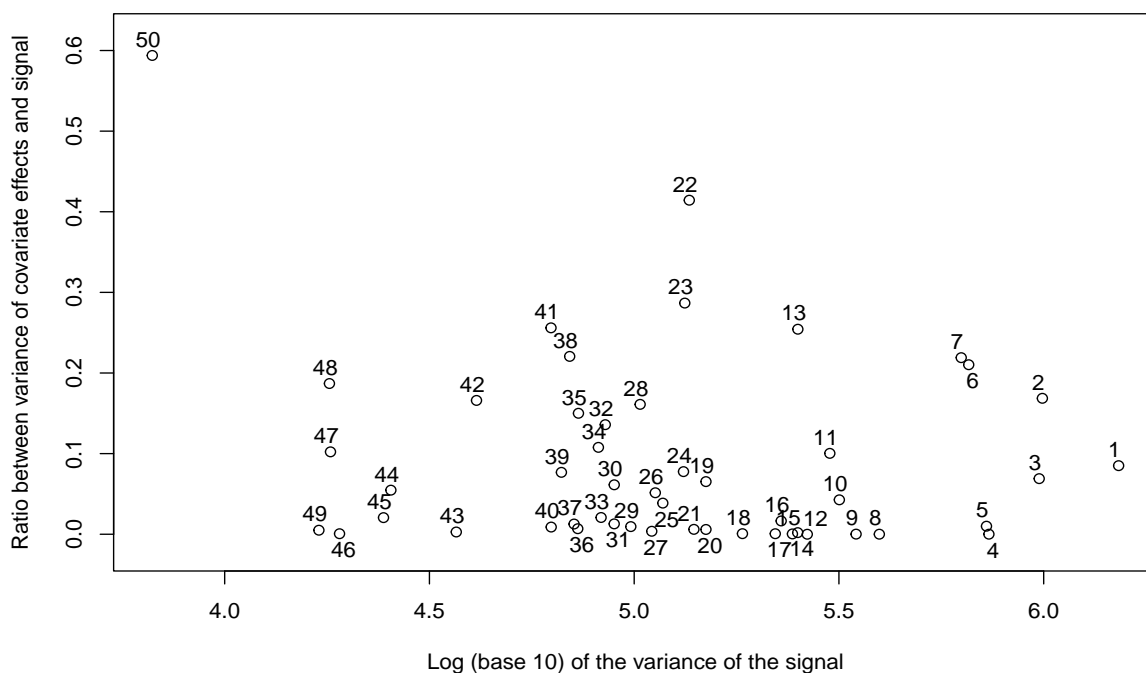
Figure 6: Scatterplot of The x-axis is the log (base 10) of the variance of the signal; the y-axis is the ratio between variance of covariate effects $\mathbf{B}_i\mathbf{Y}^{(i)}$ and signal. The number near each dot is the module identification as in $\mathbf{C}_S$.

Besides the individual and partially shared shared structures, we are able to detect global shared effects as well. Module 12 has little covariate effects, which implies that all 50 candidate mutations do not have significant global effects on 29 cancer types. However, based on Figure 7, we are able to observe several clusters that share across all the cancer types (horizontally).
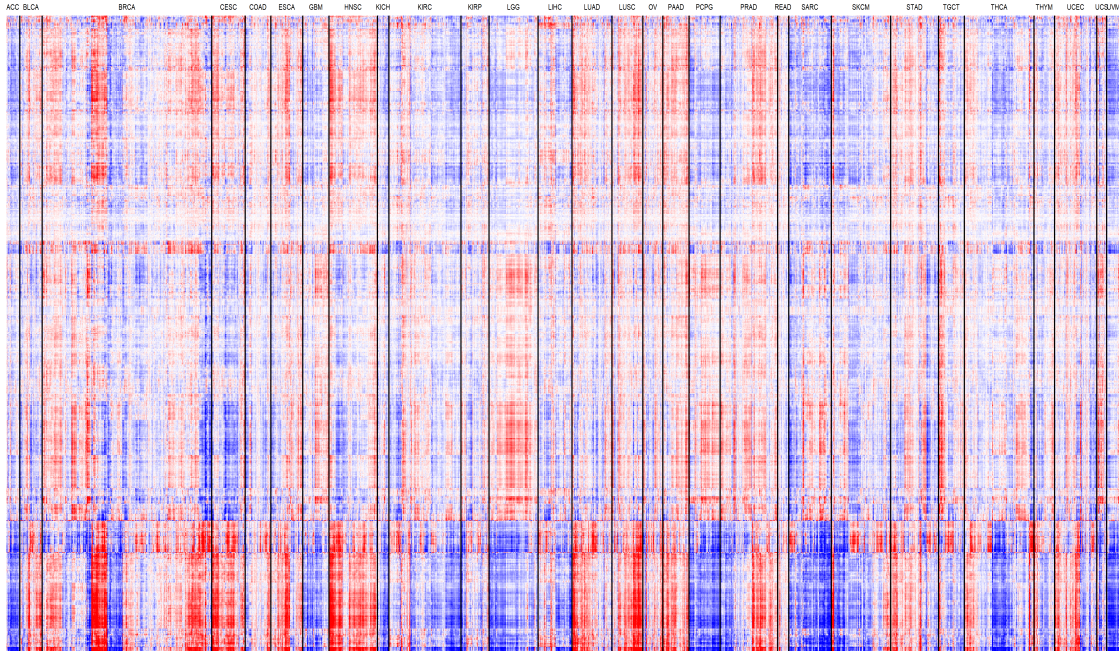
Figure 7: Heatmap for Module 12 auxiliary structure. Columns represent individual samples, while rows represent distinct gene expressions. This column-row representation is consistent across all other heatmaps in the online application results spreadsheet. Extreme values outside of 3 standard deviations are set to be threshold values. The graphics is based on the relative scale. Red colors represent high gene expressions and blue colors represents low gene expressions.

# References

Bunea, F., She, Y., and Wegkamp, M. H. (2011). Optimal selection of reduced rank estimators of high-dimensional matrices. The Annals of Statistics **39,** 1282–1309.

Cai, J.-F., Candès, E. J., and Shen, Z. (2010). A singular value thresholding algorithm for matrix completion. SIAM Journal on optimization **20,** 1956–1982.

Cancer Genome Atlas Research Network (2012). Comprehensive molecular portraits of human breast tumours. Nature **490,** 61–70.

Chen, K., Dong, H., and Chan, K.-S. (2013). Reduced rank regression via adaptive nuclear norm penalization. Biometrika **100,** 901–920.

Dolezal, J. M., Trzcinska, A., Liao, C.-Y., Kochanny, S., Blair, E., Agrawal, N., Keutgen, X. M., Angelos, P., Cipriani, N. A., and Pearson, A. T. (2021). Deep learning prediction of braf-ras gene expression signature identifies noninvasive follicular thyroid neoplasms with papillary-like nuclear features. Modern Pathology **34,** 862–874.

Feng, Q., Jiang, M., Hannig, J., and Marron, J. (2018). Angle-based joint and individual variation explained. Journal of multivariate analysis **166,** 241–265.

Gavish, M. and Donoho, D. L. (2017). Optimal shrinkage of singular values. IEEE Transactions on Information Theory **63,** 2137–2152.

Gaynanova, I. and Li, G. (2019). Structural learning and integrative decomposition of multi-view data. Biometrics **75,** 1121–1132.

Ham, S. W., Jeon, H.-Y., Jin, X., Kim, E.-J., Kim, J.-K., Shin, Y. J., Lee, Y., Kim, S. H., Lee, S. Y., Seo, S., et al. (2019). Tp53 gain-of-function mutation promotes inflammation in glioblastoma. Cell Death & Differentiation **26,** 409–425.

Hoadley, K. A., Yau, C., Hinoue, T., Wolf, D. M., Lazar, A. J., Drill, E., Shen, R., Taylor, A. M., Cherniack, A. D., Thorsson, V., et al. (2018). Cell-of-origin patterns dominate the molecular classification of 10,000 tumors from 33 types of cancer. Cell **173,** 291–304.

Hutter, C. and Zenklusen, J. C. (2018). The cancer genome atlas: creating lasting value beyond its data. Cell **173,** 283–285.

Izenman, A. J. (1975). Reduced-rank regression for the multivariate linear model. Journal of multivariate analysis **5,** 248–264.

Li, G., Liu, X., and Chen, K. (2019). Integrative multi-view regression: Bridging group-sparse and low-rank models. Biometrics **75,** 593–602.

Li, G., Yang, D., Nobel, A. B., and Shen, H. (2016). Supervised singular value decomposition and its asymptotic properties. Journal of Multivariate Analysis **146,** 7–17.

Lock, E. F., Park, J. Y., and Hoadley, K. A. (2022). Bidimensional linked matrix factorization for pan-omics pan-cancer analysis. The annals of applied statistics **16,** 193.

Mazumder, R., Hastie, T., and Tibshirani, R. (2010). Spectral regularization algorithms for learning large incomplete matrices. The Journal of Machine Learning Research **11,** 2287–2322.

Olivier, M., Hollstein, M., and Hainaut, P. (2010). Tp53 mutations in human cancers: origins, consequences, and clinical use. Cold Spring Harbor perspectives in biology **2,** a001008.

Rudelson, M. and Vershynin, R. (2010). Non-asymptotic theory of random matrices: extreme singular values. In Proceedings of the ICM 2010, pages 1576–1602. World Scientific.

Shabalin, A. A. and Nobel, A. B. (2013). Reconstruction of a low-rank matrix in the presence of gaussian noise. Journal of Multivariate Analysis **118,** 67–76.

Troyanskaya, O., Cantor, M., Sherlock, G., Brown, P., Hastie, T., Tibshirani, R., Botstein, D., and Altman, R. B. (2001). Missing value estimation methods for dna microarrays. Bioinformatics **17,** 520–525.

Wang, J. and Safo, S. E. (2021). Deep ida: A deep learning method for integrative discriminant analysis of multi-view data with feature ranking–an application to covid-19 severity. ArXiv page 2111.09964.

Yuan, M., Ekici, A., Lu, Z., and Monteiro, R. (2007). Dimension reduction and coefficient estimation in multivariate linear regression. Journal of the Royal Statistical Society: Series B (Statistical Methodology) **69,** 329–346.

Zhang, Y. and Gaynanova, I. (2022). Joint association and classification analysis of multi-view data. Biometrics **78,** 1614–1625.