

Angle-Based Joint and Individual Variation Explained

Qing Feng, Jan Hannig, Meilei Jiang*, J. S. Marron

Department of Statistics and Operations Research, University of North Carolina at Chapel Hill, Chapel Hill, NC 27599, United States

Abstract

Integrative analysis of disparate data blocks measured on a common set of experimental subjects is a major challenge in modern data analysis. This data structure naturally motivates the simultaneous exploration of the joint and individual variation within each data block resulting in new insights. For instance, there is a strong desire to integrate the multiple genomic data sets in The Cancer Genome Atlas to characterize the common and also the unique aspects of cancer genetics and cell biology for each source. In this paper we introduce Angle-Based Joint and Individual Variation Explained capturing both joint and individual variation within each data block. This is a major improvement over earlier approaches to this challenge in terms of a new conceptual understanding, much better adaption to data heterogeneity and a fast linear algebra computation. Important mathematical contributions are the use of score subspaces as the principal descriptors of variation structure and the use of perturbation theory as the guide for variation segmentation. This leads to an exploratory data analysis method which is insensitive to the heterogeneity among data blocks and does not require separate normalization. An application to cancer data reveals different behaviors of each type of signal in characterizing tumor subtypes. An application to a mortality data set reveals interesting historical lessons.

Keywords: Data integration, Heterogeneity, Perturbation theory, Principal angle, Singular value decomposition

1. Introduction

A major challenge in modern data analysis is data integration, combining diverse information from disparate data sets measured on a common set of experimental subjects. Simultaneous variation decomposition has been useful in many practical applications. For example, Kühnle (2011), Lock and Dunson (2013), Mo et al. (2013) performed integrative clustering on multiple sources to reveal novel and consistent cancer subtypes based on understanding of joint and individual variation. The Cancer Genome

*Corresponding author.

Email addresses: qing.feng1014@gmail.com (Qing Feng), jan.hannig@unc.edu (Jan Hannig), jiangm@live.unc.edu (Meilei Jiang), marron@unc.edu (J. S. Marron)

Atlas (TCGA) (Network et al., 2012) provides a prototypical example for this problem. TCGA contains disparate data types generated from high-throughput technologies. In-
 10 tegration of these is fundamental for studying cancer on a molecular level. Other types of application include analysis of multi-source metabolomic data (Kuligowski et al., 2015), extraction of commuting patterns in railway networks (Jere et al., 2014), recog-
 nition of braincomputer interface (Zhang et al., 2015), etc.

A unified and insightful understanding of the set of data blocks is expected from
 15 simultaneously exploring the joint variation representing the inter-block associations and the individual variation specific to each block. Lock et al. (2013) formulated this challenge into a matrix decomposition problem. Each data block is decomposed into three matrices modeling different types of variation, including a low-rank approxima-
 tion of the joint variation across the blocks, low-rank approximations of the individual
 20 variation for each data block, and residual noise. Definitions and constraints were pro-
 posed for the joint and individual variation together with a method named *JIVE*, see <https://genome.unc.edu/jive/> and O’Connell and Lock (2016) for Matlab and R implementations of *JIVE* respectively.

JIVE was a promising framework for studying multiple data matrices. However,
 25 the concepts of joint and individual variation were neither fully understood nor well defined. That led to problems in computation. The Lock et al. (2013) algorithm and its implementation was iterative (thus slow) and had no guarantee of achieving a so-
 lution that satisfied the definitions of *JIVE*. Another drawback of that approach was a need for arbitrary normalization of the data sets which can be hard to choose in
 30 some complicated contexts. The example in Figure B.12 in Appendix B shows this can be a serious issue. An important related algorithm named COBE was developed by Zhou et al. (2016). COBE considers a *JIVE* type decomposition as a quadratic op-
 timization problem with restrictions to ensure identifiability. While COBE removed many of the shortcomings of the original *JIVE*, it was still iterative and often required
 35 longer computation time than the Lock et al. (2013) algorithm. Neither Zhou et al. (2016) nor Lock et al. (2013) provided any theoretical basis for selection of a thresh-
 olding parameter used for separation of the joint and individual components.

A novel solution, *Angle-based Joint and Individual Variation Explained (AJIVE)*,
 is proposed here for addressing this matrix decomposition problem. It provides an
 40 efficient *angle-based algorithm* ensuring an identifiable decomposition and also an in-
 sightful new interpretation of extracted variation structure. The key insight is the use of row spaces, i.e., a focus on scores, as the principal descriptor of the joint and individual
 variation, assuming columns are the n data objects, e.g., vectors of measurements on
 patients. This focuses the methodology on variation patterns across data objects, which
 45 gives straightforward definitions of the components and thus provides identifiability.
 These variation patterns are captured by the *score subspaces* of \mathbb{R}^n . Segmentation of
 joint and individual variation is based on studying the relationship between these score
 subspaces and using perturbation theory to quantify noise effects (Stewart and Sun,
 1990).

50 The main idea of *AJIVE* is illustrated in the flowchart of Figure 1. *AJIVE* works
 in three steps. First we find a low rank approximation of each data block (shown as the
 far left color blocks in the flowchart) using SVD with e.g. a threshold selected using a
 scree plot. This is depicted (using blocks with colored dashed line boundaries) on the

left side of Figure 1 with the black arrows signifying thresholded SVD. Next, in the middle of the figure, SVD of the concatenated bases of row spaces from the first step (the gray blocks with colored boundaries) gives a joint row space (the gray box next to the circle), using a mathematically rigorous threshold derived using perturbation theory in Section 2.3. This SVD is a natural extension of Principal Angle Analysis, which is also closely related to the multi-block extension of Canonical Correlation Analysis (Nielsen, 2002) as well as to the flag means of the row spaces (Draper et al., 2014), see Section 5.2 for details. Finally, the joint and individual space approximations are found using projection of the joint row space and its orthogonal complements on the data blocks as shown as colored boundary gray squares on the right with the three joint components at the top and the individual components at the bottom.

Using score subspaces to describe variation contained in a matrix not only empowers the interpretation of analysis but also improves understanding of the problem and the efficiency of the algorithm. An identifiable decomposition can now be obtained with all definitions and constraints satisfied even in situations when individual spaces are somewhat correlated. Moreover, the need to select a tuning parameter used to distinguish joint and individual variation is eliminated based on theoretical justification using perturbation theory. A consequence is an algorithm which uses a fast built in singular value decomposition to replace lengthy iterative algorithms. For the example in Section 1.1, implemented in Matlab, the computational time of AJIVE (10.8 seconds) is about 11 times faster than the old JIVE (121 seconds) and 39 times faster than COBE (422 seconds). The computational advantages of AJIVE get even more pronounced on data sets with higher dimensionality and more complex heterogeneity such as the TCGA data analyzed in Section 4.1. For a very successful application of AJIVE on integrating fMRI imaging and behavioral data see Yu et al. (2017).

Other methods that aim to study joint variation patterns and/or individual variation patterns have also been developed. Westerhuis et al. (1998) discusses two types of methods. One main type extends traditional Principal Component Analysis (PCA), including Consensus PCA and Hierarchical PCA first introduced by Wold et al. (1987, 1996). An overview of extended PCA methods is discussed in Smilde et al. (2003). Abdi et al. (2013) discuss a multiple block extension of PCA called multiple factor analysis. This type of method computes the block scores, block loadings, global loadings and global scores.

The other main type of method is extensions of Partial Least Squares (PLS) (Wold, 1985) or Canonical Correlation Analysis (CCA) (Hotelling, 1936) that seek associated patterns between the two data blocks by maximizing covariance/correlation. For example, Wold et al. (1996) introduced multi-block PLS and hierarchical PLS (HPLS) and Trygg and Wold (2003) proposed *O2-PLS* to better reconstruct joint signals by removing structured individual variation. A multi-block extension can be found in Löfstedt et al. (2013).

Yang and Michailidis (2015) provide a very nice integrative joint and individual component analysis based on non-negative matrix factorization. Ray et al. (2014) do integrative analysis using factorial models in the Bayesian setting. Schouteden et al. (2013, 2014) propose a method called DISCO-SCA that is a low-rank approximation with rotation to sparsity of the concatenated data matrices.

A connection between extended PCA and extended PLS methods is discussed

AJIVE Path Diagram

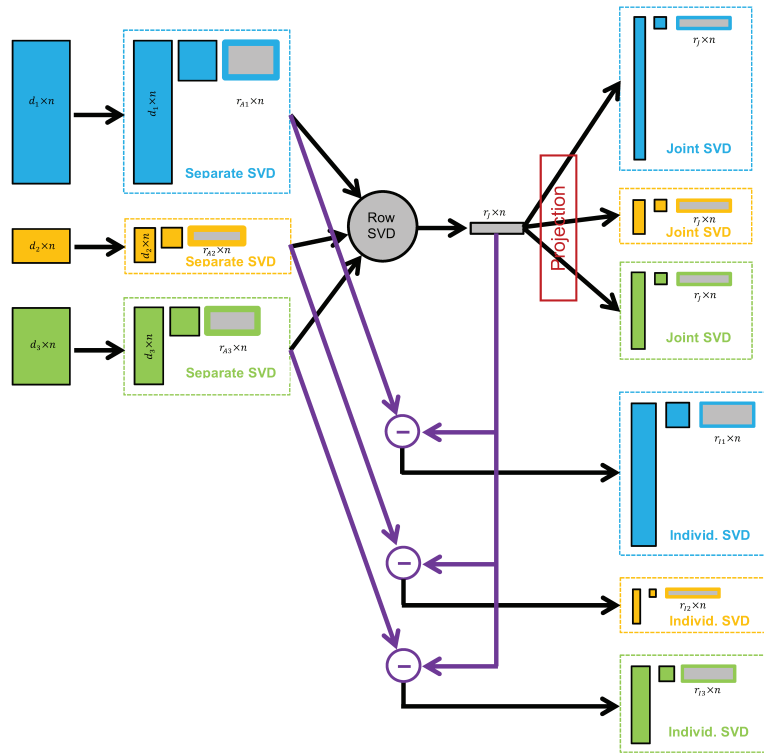


Figure 1: Flow chart demonstrating the main steps of AJIVE. First low rank approximation of each data block is obtained on the right. Then in the middle joint structure between the low rank approximations is extracted using SVD of the stacked row basis matrices. Finally, on the right, the joint components (upper) are obtained by projection of each data block onto the joint basis (middle) and the individual components (lower) come from orthonormal basis subtraction.

100 in Hanafi et al. (2011). Both types of methods provide an integrative analysis by taking the inter-block associations into account. These papers recommend use of normalization to address potential scale heterogeneity, including normalizing by the Frobenius norm, or the largest singular value of each data block etc. However, there are no consistent criteria for normalization and some of these methods have convergence problems.

105 An important point is that none of these approaches provide simultaneous decomposition highlighting joint and individual modes of variation with the goal of contrasting these to reveal new insights.

1.1. Toy Example

We give a toy example to provide a clear view of multiple challenges brought by potentially very disparate data blocks. This toy example has two data blocks, X (100 \times 100) and Y (10000 \times 100), with patterns corresponding to joint and individual structures. Such data set sizes are reasonable in modern genetic studies, as seen in Section 4.1. Figure 2 shows colormap views of these matrices, with the value of each matrix entry colored according to the color bar at the bottom of each subplot. The data has been simulated so expected row means are 0. Therefore mean centering is not necessary in this case. A careful look at the color bar scalings shows the values are almost 4 orders of magnitude larger for the top matrices. Each column of these matrices is regarded as a common data object and each row is considered as one feature. The number of features is also very different as labeled in the y-axis. Each of the two raw data matrices, X and Y in the left panel of Figure 2, is the sum of joint, individual and noise components shown in the other panels.

The joint variation for both blocks, second column of panels, presents a contrast between the left and right halves of the data matrix, thus having the same rank one score subspace. If for example the left half columns were male and right half were female, this joint variation component can be interpreted as a contrast of gender groups which exists in both data blocks for those features where color appears.

The X individual variation, third column of panels, partitions the columns into two other groups of size 50 that are arranged so the row space is orthogonal to that of the joint score subspace. The individual signal for Y contains two variation components, each driven by half of the features. The first component, displayed in the first 5000 rows, partitions the columns into three groups. The other component is driven by the bottom half of the features and partitions the columns into two groups, both with row spaces orthogonal to the joint. Note that these two individual score subspaces for X and Y are different but not orthogonal. The smallest angle between the individual subspaces is 48°.

This example presents several challenging aspects, which also appear in real data sets such as TCGA, as studied in Section 4.1. One is that the values of the features are orders of magnitude different between X and Y . There are two standard approaches to handle this, both having drawbacks. Feature by feature normalization loses information in X because Y has so many more features. Total power normalization tends to underweight the signal in Y because each feature then receives too little weight. Another important challenge is that because the individual spaces are not orthogonal, the individual signals are correlated. Correctly handling this is a major improvement of AJIVE over earlier methods.

145 The noise matrices, the right panels of Figure 2, are standard Gaussian random matrices (scaled by 5000 for X) which generates a noisy context for both data blocks and thus a challenge for analysis, as shown in the left panels of Figure 2.

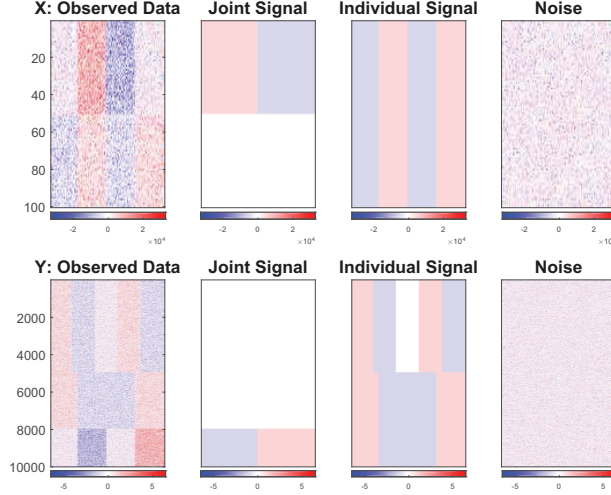


Figure 2: Data blocks X (top) and Y (bottom) in the toy example. The left panels present the observed data matrices which are a sum of the signal and noise matrices depicted in the remaining panels. Scale is indicated by color bars at the bottom of each sub-plot. These structures are challenging to capture using conventional methods due to very different orders of magnitude and numbers of features.

Simply concatenating X and Y on columns and performing a singular value decomposition on this concatenated matrix completely fails to give a meaningful joint analysis. PLS and CCA might be used to address the magnitude difference in this example and capture the signal components. However, they target common relationships between two data matrices and therefore are not able to simultaneously extract and distinguish the two types of variation. Moreover, because of its sensitivity to the strength of the signal PLS misclassifies correlated individual components as joint components. The original JIVE of Lock et al. (2013) also fails on this toy example. Details on all of these can be found in Appendix B.

The left panel of Figure 3 shows the AJIVE approximation of each data block which well captures the signal variations within both X and Y . What's more, our method correctly distinguishes the types of variation showing its robustness against heterogeneity across data blocks and correlation between individual data blocks. The approximations of both joint and individual signal are depicted in the remaining panels.

The rest of this paper is organized as follows. Section 2 describes the population model and mathematical details of the estimation approach. Results of application to a TCGA breast cancer data set and a mortality data set are presented in Section 4. Relationships between the proposed AJIVE and other methods from an optimization point of view are discussed in Section 5.

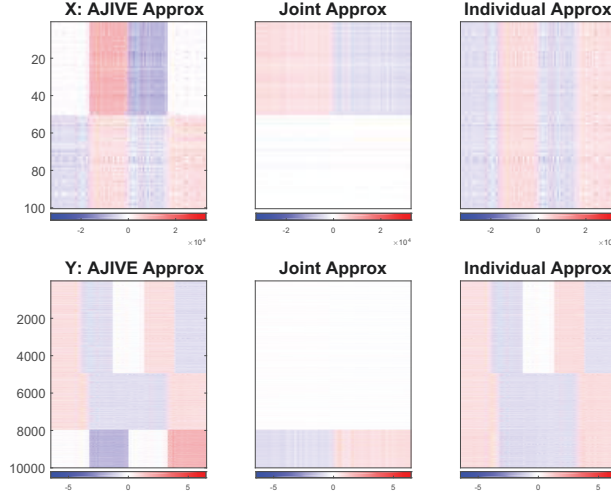


Figure 3: AJIVE approximation of the data blocks X and Y in the toy example are shown in the first column, with the joint and individual signal matrices depicted in the remaining columns. Both quite diverse types of variations are well captured for each data block by AJIVE, in contrast to other usual methods as seen in Appendix B.

2. Proposed Method

In this section the details of the new proposed AJIVE are discussed. The population model is proposed in Sections 2.1 and 2.2. The theoretical foundations based on matrix perturbation theory from linear algebra (Stewart and Sun, 1990) are given in Section 2.3. These theoretical results motivate our estimation approach which is proposed in Section 2.4.

2.1. Population Model - Signal

Matrices $\{X_k, k = 1, \dots, K\}$ ($d_k \times n$) are a set of data blocks for study, e.g. the colored blocks on the left of Figure 1. The columns are regarded as data objects, one vector of measurements for each experimental subject, while rows are considered as features. All X_k therefore have the same number of columns and perhaps a different number of rows.

Each X_k is modeled as low rank true underlying signals A_k perturbed by additive noise matrices E_k . Each low rank signal A_k is the sum of two matrices containing joint and individual variation, denoted as J_k and I_k respectively for each block

$$\begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_K \end{bmatrix} = \begin{bmatrix} A_1 \\ A_2 \\ \vdots \\ A_K \end{bmatrix} + \begin{bmatrix} E_1 \\ E_2 \\ \vdots \\ E_K \end{bmatrix} = \begin{bmatrix} J_1 \\ J_2 \\ \vdots \\ J_K \end{bmatrix} + \begin{bmatrix} I_1 \\ I_2 \\ \vdots \\ I_K \end{bmatrix} + \begin{bmatrix} E_1 \\ E_2 \\ \vdots \\ E_K \end{bmatrix}. \quad (1)$$

Our approach focuses on the vectors in the row space of our matrices. In this context these vectors are often called *score vectors* and the row space of the matrix is often

called *score subspace* ($\subset \mathbb{R}^n$). Therefore, the row spaces of the matrices capturing joint variation, i.e., joint matrices, are defined as sharing a common score subspace denoted as $\text{row}(J)$

$$\text{row}(J_k) = \text{row}(J), \quad k = 1, \dots, K.$$

The individual matrices are individual in the sense that the intersection of their score subspaces is the zero vector space, i.e.

$$\bigcap_{k=1}^K \text{row}(I_k) = \{\vec{0}\}, \quad k = 1, \dots, K.$$

This means that there is no non-trivial common row pattern in every individual score subspaces across blocks.

To ensure an identifiable variation decomposition we assume $\text{row}(J) \subset \text{row}(A_k)$, which also implies $\text{row}(I_k) \subset \text{row}(A_k)$, for $k = 1, \dots, K$. Orthogonality between the score subspaces of matrices containing joint and individual variation is also assumed. In particular, $\text{row}(J) \perp \text{row}(I_k)$, $k = 1, \dots, K$. Note that orthogonality between individual matrices $\{I_k, k = 1, \dots, K\}$ is *not* assumed as it is not required for the model to be uniquely determined.

Under these assumptions, the model is identifiable in the sense:

Lemma 2.1. *Given a set of matrices $\{A_k, k = 1, \dots, K\}$, there are unique sets of matrices $\{J_k, k = 1, \dots, K\}$, and $\{I_k, k = 1, \dots, K\}$ so that:*

1. $A_k = J_k + I_k, k = 1, \dots, K$
2. $\text{row}(J_k) = \text{row}(J) \subset \text{row}(A_k), k = 1, \dots, K$
3. $\text{row}(J) \perp \text{row}(I_k), k = 1, \dots, K$
4. $\bigcap_{k=1}^K \text{row}(I_k) = \{\vec{0}\}.$

The proof is provided in Appendix A. This model has enhanced the matrix decomposition idea proposed in Lock et al. (2013) by providing a clearer mathematical framework and precise understanding of the different types of variation. In particular, Lock et al. (2013) imposed rank constraints on the joint matrices i.e. $\text{rank}(J_k)$ are the same for all data blocks but did not clearly stipulate the need of a common score space. An unnecessary orthogonality among individual matrices was further suggested, although not explicitly enforced in the estimation, for ensuring a well defined decomposition.

2.2. Population Model - Noise

The additive noise matrices are assumed to follow an isotropic error model where the energy of projection is invariant to direction in both row and column spaces. Important examples include the multivariate standard normal distribution and the multivariate student t-distribution (Kotz and Nadarajah, 2004). The singular values of each noise matrix are assumed to be smaller than the smallest singular values of each signal to give identifiability.

The assumption on the noise distribution here is less strong than the classical i.i.d. Gaussian random matrix, and only comes into play when determining the number of joint components. The estimation approach given in Section 2.3 reconstructs each signal matrix based on SVD and bootstrap-like resampling and thus is relatively insensitive to the error distribution.

2.3. Theoretical Foundations

The main challenge is segmentation of the joint and individual variation in the presence of noise which individually perturbs each signal. Let $\{\tilde{A}_k, k = 1, \dots, K\}$ be noisy approximations of $\{A_k, k = 1, \dots, K\}$ respectively. The subspaces of joint variation within the approximations \tilde{A}_k , while expected to be similar, are no longer exactly the same due to noise. If some subspaces of $\{\tilde{A}_k, k = 1, \dots, K\}$ are very close, they can be considered as estimates of the common score subspace under different perturbations. Application of the results of the *Generalized sin θ Theorem* (Wedin, 1972) is proposed to decide when a set of subspaces are close enough to be regarded as estimates of the joint score space. Based on this theorem, the number of joint components can be determined resulting in an appropriate segmentation.

Take the approximation \tilde{A}_k of A_k as an example of perturbation of the score space of each matrix. For consistency with the Generalized sin θ Theorem, we use the following pseudometric as a notion of distance between theoretical and perturbed subspaces. Let $\mathcal{Q}_k, \tilde{\mathcal{Q}}_k$ be the l dimensional score subspaces of \mathbb{R}^n respectively for the matrix A_k and its approximation \tilde{A}_k . The corresponding symmetric projection matrices are $P_{\mathcal{Q}_k}$ and $P_{\tilde{\mathcal{Q}}_k}$. The distance between the two subspaces is defined as the difference of the projection matrices under the operator L^2 norm, i.e., $\rho(\mathcal{Q}_k, \tilde{\mathcal{Q}}_k) = \|P_{\mathcal{Q}_k} - P_{\tilde{\mathcal{Q}}_k}\|$ (Stewart and Sun, 1990).

An insightful understanding of this pseudometric $\rho(\mathcal{Q}_k, \tilde{\mathcal{Q}}_k)$ comes from a principal angle analysis (Jordan, 1875; Hotelling, 1936) of the subspaces \mathcal{Q}_k and $\tilde{\mathcal{Q}}_k$. Denote the principal angles between \mathcal{Q}_k and $\tilde{\mathcal{Q}}_k$ as $\Theta(\mathcal{Q}_k, \tilde{\mathcal{Q}}_k) = \{\theta_{k,1}, \dots, \theta_{k,l}\}$ with $\theta_{k,1} \geq \theta_{k,2} \geq \dots \geq \theta_{k,l}$. The pseudometric $\rho(\mathcal{Q}_k, \tilde{\mathcal{Q}}_k)$ is equal to the sine of the maximal principal angle, i.e., $\sin \theta_{k,1}$. This suggests that the largest principal angle between two subspaces can indicate their closeness, i.e. distance. Under a slight perturbation, the largest principal angle between \mathcal{Q}_k , a theoretical subspace, and $\tilde{\mathcal{Q}}_k$, its perturbed subspace, is expected to be small.

The pseudometric $\rho(\mathcal{Q}_k, \tilde{\mathcal{Q}}_k)$ can be also written as

$$\rho(\mathcal{Q}_k, \tilde{\mathcal{Q}}_k) = \|(I - P_{\mathcal{Q}_k})P_{\tilde{\mathcal{Q}}_k}\| = \|(I - P_{\tilde{\mathcal{Q}}_k})P_{\mathcal{Q}_k}\|$$

which brings another useful understanding of this definition. It measures the relative deviation of the signal variation from the theoretical subspace. Accordingly, the similarity/closeness between the subspaces and its perturbation can be written as $\|P_{\mathcal{Q}_k}P_{\tilde{\mathcal{Q}}_k}\|$ and is equal to the cosine of the maximal principal angle defined above, i.e. $\cos \theta_{k,1}$. Hence, $\sin^2 \theta_{k,1}$ indicates the percentage of signal deviation and $\cos^2 \theta_{k,1}$ tells the percentage of remaining signal in the theoretical subspace.

The generalized sin θ theorem provides a bound for the distance between a subspace and its perturbation, e.g., the subspaces \mathcal{Q}_k and $\tilde{\mathcal{Q}}_k$. This bound quantifies how

the theoretical subspace \mathcal{Q}_k is affected by noise. In particular, the following is an
 250 adaptation of the the generalized $\sin \theta$ theorem to our setup:

Theorem 2.2 (Wedin, 1972). *For $k = 1, \dots, K$, signal matrix A_k is perturbed by additive noise E_k . Let $\theta_{k,1}$ be the largest principal angle for the subspace of signal A_k and its approximation \tilde{A}_k . Denote the SVD of \tilde{A}_k as $\tilde{U}_k \tilde{\Sigma}_k \tilde{V}_k^T$. The distance between the subspaces of A_k and \tilde{A}_k , $\rho(\mathcal{Q}_k, \tilde{\mathcal{Q}}_k)$ i.e. sines of $\theta_{k,1}$, is bounded*

$$\rho(\mathcal{Q}_k, \tilde{\mathcal{Q}}_k) = \sin \theta_{k,1} \leq \frac{\max(\|E_k \tilde{V}_k\|, \|E_k^T \tilde{U}_k\|)}{\sigma_{\min}(\tilde{\Sigma}_k)}, \quad (2)$$

where $\sigma_{\min}(\tilde{\Sigma}_k)$ is the smallest singular value of \tilde{A}_k .

This bound measures how far the perturbed space can be away from the theoretical one. The deviation is bounded by the maximal value of noise energy on the column and row spaces and also the smallest signal singular values. This is consistent with
 255 the intuition that a deviation distance, i.e., a largest principal angle, is small when the signal is strong and perturbations are weak.

Notice that the bound in Theorem 2.2 is applicable but cannot be directly used for data analysis since the error matrices E_k are not observable. As the error matrices are assumed to be isotropic, we propose to re-sample noisy directions from the residuals
 260 of the low rank approximations. The operator L^2 norm of the error related terms in Theorem 2.2 can thus be estimated by projecting the observed data onto the subspace spanned by re-sampled directions. This re-sampling based method can also provide prediction intervals for these perturbation bounds. More details of estimating the perturbation bound will be discussed next.

265 2.4. Estimation Approach

A three-step algorithm as illustrated in Figure 1 is outlined below. This algorithm uses SVD in each step. As a basic illustration for each step we use the toy example described in Section 1. Details for each step appear in the following subsections.

Step 1: Signal Space Initial Extraction: Even though the signal components
 270 $\{A_k, k = 1, \dots, K\}$ are low rank, the data matrices $\{X_k, k = 1, \dots, K\}$ are usually of full rank due to presence of noise. SVD works as a signal extraction device in this step, keeping components with singular values greater than selected thresholds individually for each data block.

When selecting these thresholds, one needs to be aware of a bias/variance like trade-off. Setting the threshold too high will provide an accurate estimation of the parts of
 275 the joint space that are included in the low rank approximation. The downside is that significant portions of the joint signal might be thresholded out. This could be viewed as a low variance high bias situation. If the threshold is set low than it is likely that the joint signal is included in all of the blocks. However, the precision of the segmentation
 280 in the next step can deteriorate to the point that most of the joint space selected is driven by noise. This can be viewed as the low bias high variance situation.

Most off the shelf automatic procedures for low rank matrix approximation have as their stated goal signal reconstruction and prediction which based on our experience

lead to thresholds that are too small. This is sensible as adding a little bit more noise usually helps prediction but it has bad effects on signal segmentation. We therefore recommend taking a multi-scale perspective and trying several threshold choices, for example, by manually finding several relatively big jumps in a scree plot. Figure 4 shows the scree plots of each data block for the toy example in Section 1. The left scree plot for X suggests a selection of rank as 2 and the right one for Y suggests rank 3, since in both cases those components stand out while the rest of the singular values decay slowly showing no clear jump.

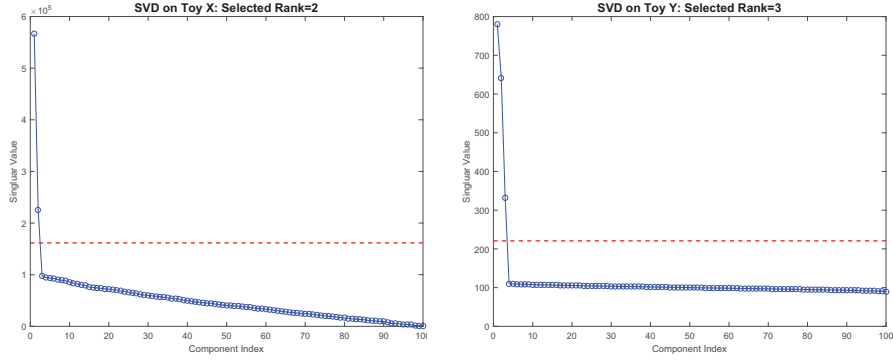


Figure 4: Scree plots for the toy data sets X (left) and Y (right). Both plots display the singular values associated with a component in descending order versus the index of the component. The components with singular values above the dashed red threshold line are regarded as the initial signal components in the first step of AJIVE.

Let $\{\tilde{r}_k, k = 1, \dots, K\}$ be the initial estimates of the signal ranks $\{r_k, k = 1, \dots, K\}$. In the toy example $\tilde{r}_1 = 2$ (for X) and $\tilde{r}_2 = 3$ (for Y). Each data block has a low rank approximation, \tilde{A}_k (represented in Figure 1 as the boxes with dashed colored boundaries on the left), which is the initial estimate of the signal matrix A_k , $k = 1, \dots, K$. The estimate is decomposed as

$$\tilde{A}_k = \tilde{U}_k \tilde{\Sigma}_k \tilde{V}_k^T, \quad (3)$$

where \tilde{U}_k contains the left singular vectors that correspond to the largest \tilde{r}_k singular values respectively for each data block. The initial estimate of the signal score space, denoted as $\text{row}(\tilde{A}_k)$, is spanned by the right singular vectors in \tilde{V}_k (shown as gray boxes with colored boundaries on the left).

Step 2A: Score Space Segmentation of Two Data Blocks: For a clear introduction to the basic idea of score space segmentation, the two-block special case ($K = 2$) is first studied. The goal is to use the low rank approximations \tilde{A}_k from equation (3) to obtain estimates of the common joint and individual score subspaces. Due to the presence of noise, the components of $\text{row}(\tilde{A}_k)$, $k = 1, 2$, corresponding to the underlying joint space, no longer are the same, but should have a relatively small angle. Similarly, the components corresponding to the underlying individual spaces are expected to have

a relatively large angle. This motivates the use of principal angle analysis to separate the joint from the individual components. The following Lemma 2.3 provides a bound
 305 on the largest allowable principal angle of the joint part of the initial estimated spaces.

Lemma 2.3. *Let ϕ be the largest principal angle between two subspaces that are each a perturbation of the common row space within $\text{row}(\tilde{A}_1)$ and $\text{row}(\tilde{A}_2)$. That angle is bounded by*

$$\sin \phi \leq \sin(\theta_{1,1} + \theta_{2,1}) \quad (4)$$

in which $\theta_{1,1}$ and $\theta_{2,1}$ are the angles given in Theorem 2.2.

The proof is provided in Appendix A. As mentioned in Section 2.3, the perturbation bounds of each $\theta_{k,1}$ require the estimation of terms $\|E_k \tilde{V}_k\|$, $\|E_k^T \tilde{U}_k\|$ for $k = 1, 2$. These terms are the measurements of energies of noise matrices projected onto the signal column and row spaces. Since an isotropic error model is assumed, the distribution
 310 of energy of the noise matrices in arbitrary directions are supposed to be equal. Thus, if we sample directions orthogonal to the estimated signal and use the observe data X_k instead of E_k , this should provide a good estimator of the *distribution* of the unobserved terms $\|E_k \tilde{V}_k\|$, $\|E_k^T \tilde{U}_k\|$. To this end, denote \tilde{V}_k^\perp ($n \times (\min(d_k, n) - \tilde{r}_k)$) and \tilde{U}_k^\perp ($d_k \times (\min(d_k, n) - \tilde{r}_k)$) as the respective orthonormal bases of the row and column subspaces of the residual matrices from the low rank approximations in equation
 315 (3) and resample non-zero column vectors from the matrices \tilde{V}_k^\perp and \tilde{U}_k^\perp . Notice that these matrices are not full rank.

Take the term $\|E_k \tilde{V}_k\|$ as an example for illustration. Given the \tilde{r}_k number of column vectors resampled without replacement from \tilde{V}_k^\perp , denoted as V^* , the observed
 320 data block X_k is projected onto the subspace spanned by V^* , written as $X_k V^*$. The distribution of the operator L^2 norm $\|X_k V^*\|$ approximates the distribution of the unknown $\|E_k \tilde{V}_k\|$. Thus we resample 1000 of $\|X_k V^*\|$ and use the quantiles to provide both a point estimate and a simulated prediction interval for $\|E_k \tilde{V}_k\|$. This can be
 325 similarly applied to $\|E_k^T \tilde{U}_k\|$ for $k = 1, 2$, resulting in a prediction interval for the perturbation bound. Typically the median is chosen as the estimate of the angle bound for exploratory analysis. This will result in 50% confidence that all joint components are included. For certain cases where finding most of the joint components is desired, the 95th percentile of these estimated terms can be used to derive a conservative angle
 330 threshold, resulting in at least 95% confidence of finding all joint components that were included in Step 1.

To investigate the validity of this approximation to Theorem 2.2 we performed a small scale simulation study based on the example of Section 1.1. We generated 10,000 independent copies of the data sets X (100×100 , true signal rank 2) and Y
 335 (1000×100 , true signal rank 3). Then for several low rank approximations (columns of Table 1) we calculated the estimate of the angle between the true signal and the low rank approximation using the approximation above. Table 1 reports the percentage of the times the corresponding quantile of the resampled estimate is bigger than the true angle for the matrix X . We see that the performance for the square matrix X is satisfactory as the empirical percentages are close to the nominal values. Corresponding
 340 empirical percentages for the high dimensional low sample size data set Y are all 100 %, and thus are not shown. This is caused by the fact that Wedin's bound can be very

Table 1: Coverages of the prediction intervals of the true angle between signal and low rank approximation for X . Rows are nominal levels. Columns are ranks of approximation (where 2 is the correct rank). The simulation shows good performance for the square matrix X .

	1	2	3
50%	91.9%	63.6%	100.0%
90%	100.0%	89.6%	100.0%
95%	100.0%	93.7%	100.0%
99%	100.0%	98.0%	100.0%

conservative if the matrix is far from square. We also compared the estimated angle between the perturbed joint spaces in X and Y with the actual angle. As expected
 345 the estimate remains conservative with all of the 10000 estimates being larger than the true angle. Recent work of Cai and Zhang (2016) may provide a potential approach for improvement.

One of the ways of computing the principal angles between $\text{row}(\tilde{A}_1)$ and $\text{row}(\tilde{A}_2)$ is to perform SVD on a concatenation of their right singular vector matrices (Miao and Ben-Israel, 1992), i.e.,

$$M \triangleq \begin{bmatrix} \tilde{V}_1^T \\ \tilde{V}_2^T \end{bmatrix} = U_M \Sigma_M V_M^T, \quad (5)$$

where the singular values, on the diagonal of Σ_M , determine the principal angles, $\Phi(\text{row}(\tilde{A}_1), \text{row}(\tilde{A}_2)) = \{\phi_1, \dots, \phi_l\}$ as

$$\phi_i = \arccos((\sigma_{M,i})^2 - 1), \quad i = 1, \dots, \min(\tilde{r}_1, \tilde{r}_2). \quad (6)$$

Given a left singular vector $U_{M,i}$ denoted as \vec{u} , a pair of principal vectors $\{\vec{p}_i, \vec{q}_i\}$ in each subspace can be constructed by projecting \tilde{V}_1 and \tilde{V}_2 onto the vector \vec{u} . Denote
 350 \vec{u} as the concatenation of $[\vec{u}_1; \vec{u}_2]$. Note that the length of \vec{u}_1 is equal to the number of columns of \tilde{V}_1 and similarly for the other part. The principal vectors in each subspace can be written as $\vec{p}_i = \tilde{V}_1 \vec{u}_1$ and $\vec{q}_i = \tilde{V}_2 \vec{u}_2$ respectively. The angle between the pair of principal vectors θ_i is equal to the principal angle computed from the singular value corresponding to \vec{u} .

355 As discussed in Section 5.2 the vector \vec{v}_i , the corresponding right singular vector of V_M , points in the same direction as the sum of principal vector pairs of each subspace. When the principal angle ϕ_i is smaller than the perturbation bound ϕ , this right singular vector is taken as an estimate of the theoretical joint direction.

This SVD decomposition can be understood as a tool that sorts pairs of directions
 360 within the two subspaces in increasing order of the angle between each pair. When the corresponding principal angle is smaller than the perturbation bound ϕ , the pair of principal vectors can be considered as noisy versions of the same joint direction. Assume there are \tilde{r}_J principal angles smaller than the bound ϕ . The first \tilde{r}_J singular vectors \vec{v}_i are used as the natural orthonormal basis of the estimated joint score subspace.

365 The left panel of Figure 5 depicts the principal angles of the concatenated right singular vector matrices for the toy example in Section 1.1. Since the initial estimates

of r_x and r_y are 2 and 3, there are only two potential components for joint variation. The associated principal angles between the initially estimated signal row spaces are labeled next to the first two components as 10.99° and 47.11° . The estimated bound
 370 on the principal angle in Lemma 2.3 based on 50% prediction is 31.29° for this toy example, shown as the central red dashed line. The 5% and 95% one-sided prediction intervals of the angle bound are $[0, 30.00]$ and $[0, 32.92]$ degrees, shown as green dashed lines. Each provides a respective 5% and 95% chance for including all the joint components. This provides a clear indication that the number of joint components
 375 should be $\tilde{r}_J = 1$. The corresponding first right singular vector of M will be taken as the joint score vector.

Step 2B: Score Space Segmentation of Multiple Data Blocks: To generalize the above idea to more than two blocks, the key is to focus more on singular values than on angles in the equation (6). In other words, instead of finding an upper bound on an
 380 angle, we will focus on a lower bound on the remaining energy as expressed by the sum of the squared singular values. Hence, an analogous SVD will be used for studying the closeness of multiple initial signal score subspace estimates.

For the vertical concatenation of right singular vector matrices

$$M \triangleq (\tilde{V}_1, \dots, \tilde{V}_K)^T = U_M \Sigma_M V_M^T. \quad (7)$$

SVD sorts the directions within these K subspaces in increasing order of amount of deviation from the theoretical joint direction. The squared singular value $\sigma_{M,i}^2$ indicates
 385 the total amount of variation explained in the common direction $V_{M,i}^T$ in the score subspace of \mathbb{R}^n . A large value of $\sigma_{M,i}^2$ (close to K) suggests that there is a set of basis vectors within each subspace that are close to each other and thus are potential noisy versions of a common joint score vector. A threshold on singular values is needed to segment the joint components. This is done in Lemma 2.4.

Lemma 2.4. *Let θ_k be the bound on the principal angles between the theoretical subspace $\text{row}(A_k)$ and its perturbation $\text{row}(\tilde{A}_k)$ for K data blocks from Theorem 2.2. The squared singular values ($\sigma_{M,i}^2$) corresponding to the estimates of joint components satisfy*

$$\sigma_{M,i}^2 \geq K - \sum_{k=1}^K \sin^2 \theta_k \geq K - \sum_{k=1}^K \left(\frac{\max(\|E_k \tilde{V}_k\|, \|E_k^T \tilde{U}_k\|)}{\sigma_{\min}(\tilde{\Sigma}_k)} \right)^2. \quad (8)$$

390 The proof is provided in Appendix A. This lower bound is independent of the variation magnitudes. This property makes AJIVE insensitive to scale heterogeneity across each block when extracting joint variation information.

As above, the terms $\|E_k \tilde{V}_k\|$, $\|E_k^T \tilde{U}_k\|$ are resampled to derive a point estimate and prediction interval for the threshold. As for the two-block case, if there were \tilde{r}_J
 395 singular values selected, the first \tilde{r}_J right singular vectors are used as the basis of the estimate of $\text{row}(J)$.

The right panel of Figure 5 depicts the first 2 singular values of the vertical concatenated matrix M for the toy example. This is an analysis of the same data, but performed on the scale of squared singular values instead of principal angles. The associated squared singular values are labeled next to these two component as 1.98 and
 400

1.68. The estimated threshold, targeting median prediction, is 1.85 for the toy example. This threshold together with its 5% and 95% one sided prediction intervals, $[1.86, +\infty]$ and $[1.84, +\infty]$ respectively, strongly suggest that the number of joint components \hat{r}_J should be 1.

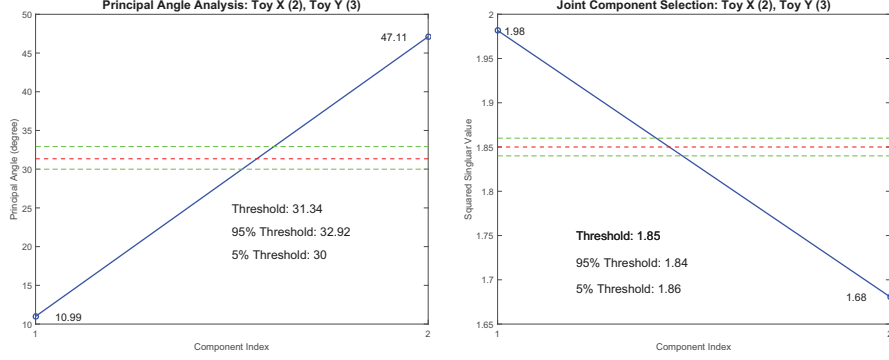


Figure 5: Left panel: Principal angles between the initial estimates of signal row spaces. The bound for the largest angle is 31.29 degree, suggesting the existence of one joint component. To indicate the uncertainty, the 5% and 95% one-sided prediction intervals of the angle threshold are also shown. Right panel: Squared singular values plot of the vertical concatenated matrix M for the toy example. Both thresholds correctly capture the underlying structure of this toy example with the selection of one joint component.

405 **Step 3: Final Decomposition:** Based on the estimate of the joint row space, matrices containing joint variation in each data block can be reconstructed by projecting X_k onto this estimated space. Define the matrix \tilde{V}_J as $[\tilde{v}_{M,1}, \dots, \tilde{v}_{M,\hat{r}_J}]$, where $\tilde{v}_{M,i}$ is the i^{th} column in the matrix V_M . To ensure that all components continue to satisfy the identifiability constraints from Section 2.2, we check that, for all the blocks, each
 410 $\|X_k \tilde{v}_{M,i}\|$ is also above the corresponding threshold used in Step 1. If the constraint is not satisfied for any block, that component is removed from \tilde{V}_J . A real example of this happens in Section 4.1. An important point is that this removal can happen even when there is a common joint structure in all but a few blocks.

Denote \hat{V}_J as the matrix \tilde{V}_J after the removal and \hat{r}_J as the final joint rank. The projection matrix onto the final estimated joint row space $\text{row}(\hat{J})$ is $P_J = \hat{V}_J$, represented as the red rectangle in Figure 1. The estimates of joint variation matrices in each block are

$$\hat{J}_k = X_k P_J, \quad k = 1, \dots, K.$$

The row space of joint structure is orthogonal to the row spaces of each individual structure. Therefore, the original data blocks are projected to the orthogonal space of $\text{row}(\hat{J})$. The projection matrix onto the orthogonal space of $\text{row}(\hat{J})$ is $P_J^\perp = I - P_J$ and the projections of each data block are denoted as X_k^\perp respectively for each block i.e.

$$X_k^\perp = X_k P_J^\perp.$$

These projections are represented as the circled minus signs in Figure 1.

415 Finally we threshold this projection by performing SVD on $\{X_k^\perp, k = 1, \dots, K\}$. The components with singular values larger than the first thresholds from Section 2.4

are kept as the individual components, denoted as $\{\hat{I}_k^\perp, k = 1, \dots, K\}$. The remaining components of each SVD are regarded as an estimate of the noise matrices.

By taking a direct sum of the estimated row spaces of each type of variation, denoted by \oplus , the estimated signal row spaces are

$$\text{row}(\hat{A}_k) = \text{row}(\hat{J}) \oplus \text{row}(\hat{I}_k)$$

with rank $\hat{r}_k = \hat{r}_J + \hat{r}_{I_k}$ respectively for $k = 1, \dots, K$.

420 Due to this adjustment of directions of the joint components, these final estimates of signal row spaces may be different from those obtained in the initial signal extraction step. Note that even the estimates of rank \hat{r}_k might also differ from the initial estimates \tilde{r}_k .

3. Post AJIVE Data Representation

Given the variation decompositions of each data block, as shown on the right side of Figure 1, several types of post AJIVE representations are available for exploring the joint and individual variation patterns. The estimates of the colored joint matrices within each data block are represented by SVD

$$\hat{J}_k = \hat{U}_J^k \hat{\Sigma}_J^k \hat{V}_J^{kT}, \quad k = 1, \dots, K \quad (9)$$

425 in which \hat{V}_J^k are the $n \times \hat{r}_J$ joint score matrices. Note that the singular values $\hat{\Sigma}_J^k$ can be completely different across k , since they are driven by the score variation pattern and can reflect very different amounts of variation between the blocks. The loading matrices \hat{U}_J^k ($d_k \times \hat{r}_J$) respectively specify distinct \hat{r}_J -dimension loading subspaces of \mathbb{R}^{d_k} for each block k .

430 There are three important matrix representations of the information in the JIVE joint output (i.e. the three boxes on the upper right, with colored dashed boundaries), with differing uses in post AJIVE analyses.

1. *Full Matrix Representation*. For applications where the original features are the main focus (such as finding driving genes) the full matrix representations \hat{J}_k ($d_k \times n$), $k = 1, \dots, K$ are most useful. Thus the JIVE output is the product of all three blocks in each dashed box. Examples are shown in Figure 3.
- 435 2. *Block Specific Score (BSS)*. For applications where the relationships between subjects are the main focus (such as discrimination between subtypes) large computational gains are available by using the much lower dimensional representations $\hat{\Sigma}_J^k \hat{V}_J^{kT}$ ($\hat{r}_J \times n$). In this case the JIVE output is the product of the right two blocks in each dashed box. This results in no loss of information when rotation invariant methods are used.
- 440 3. *Common Normalized Score (CNS)*. When it is desirable to study the component of joint behavior that is separate from the within block variation (such as evaluating the relationship between data objects), the analysis should focus on a common basis of $\text{row}(\hat{J})$, namely \hat{V}_J^T ($\hat{r}_J \times n$) from Section 2.4. Hence, the JIVE output is only the block shown with gray interior.
- 445

The relationship between BSS and CNS is analogous to that of the traditional covariance (i.e PLS) and correlation (i.e CCA) modes of analysis.

450 Furthermore, different representations provide different views of the loadings. The full matrix representation and BSS naturally obtain the information from the loading matrix \hat{U}_J^k . CNS gives a different representation of the loadings. Given the common basis of $\text{row}(\hat{J})$, one can perform regression for \hat{J}_k on each score vector in \hat{V}_J , from which the standardized coefficient vector can be taken as the CNS loading. By doing
455 this, there is no guarantee of orthogonality between CNS loading vectors. However, the loadings are linked across blocks by their common scores. Therefore, in this CNS case, the standardized regression coefficients are recommended for use instead of the classical loadings.

The individual variation within blocks can be similarly analyzed resulting in both BSS and CNS analyses for the individual components as indicated in the lower right two blocks in Figure 1. When original features are important, the full matrix

$$\hat{I}_k = \hat{U}_I^k \hat{\Sigma}_I^k \hat{V}_I^{kT}, \quad k = 1, \dots, K$$

with dimension $d_k \times n$ are available. Otherwise large computational savings are available from the BSS version $\hat{\Sigma}_I^k \hat{V}_I^{kT}$ ($\hat{r}_{I_k} \times n$), $k = 1, \dots, K$. For studying scale free behaviors, use the *Individual Normalized Score (INS)* \hat{V}_I^{kT} ($\hat{r}_{I_k} \times n$). For individual components, the matrix \hat{U}_I^k can be taken as loadings for all three representations as the INS matrices cannot be the same.

4. Data Analysis

465 In this section, we apply AJIVE to two real data sets, TCGA breast cancer in Sections 4.1 and Spanish mortality in Section 4.2.

4.1. TCGA Data

A prominent goal of modern cancer research, of which The Cancer Genome Atlas (Network et al., 2012) is a major resource, is the combination of biological insights
470 from multiple types of measurements made on common subjects.

TCGA provides prototypical data sets for the application of AJIVE. Here we study the 616 breast cancer tumor samples from Ciriello et al. (2015), which had a common measurement set. For each tumor sample, there are measurements of 16615 gene expression features (GE), 24174 copy number variations features (CN), 187 reverse
475 phase protein array features (RPPA) and 18256 mutation features (Mutation). These data sources have very different dimensions and scalings.

The tumor samples are classified into four molecular subtypes: Basal-like, HER2, Luminal A and Luminal B. An integrative analysis targets the association among the features of these four disparate data sources that jointly quantify the differences between tumor subtypes. In addition, identification of driving features for each source
480 and subtype is obtained from studying loadings.

In the first step of AJIVE we selected low rank approximations of dimensions 11 (GE), 6 (CN), 8 (RPPA) and 12 (Mutation). This gave us the most interpretable and

insightful analysis resulting in one joint component. After selection of the threshold in step 1, it took AJIVE 400 seconds (6.7 minutes) to finish steps 2 and 3.

In the second AJIVE step, the one sided 95% prediction interval suggested selection of two joint components. However, the third step indicated dropping one joint component, because the norm of the projection of the mutation data on that direction, i.e. the second CNS, is below the threshold from Step 1. This result of one joint component was consistent with the expectation of cancer researchers, who believe the mutation component has only one interesting mode of variation. A careful study of all such projections shows that the other data types, i.e. GE, CN and RPPA, do have a common second joint component as discussed at the end of this section. The association between the CNS and genetic subtype differences is visualized in the left panel of Figure 6. The dots are a jitter plot of the patients, using colors and symbols to distinguish the subtypes (Blue for Basal-like, cyan for HER2, red for Luminal A and magenta for Luminal B). Each symbol is a data point whose horizontal coordinate is the value and vertical coordinate is the height based on data ordering. The curves are Gaussian kernel density estimates i.e. smoothed histograms, which show the distribution of the subtypes.

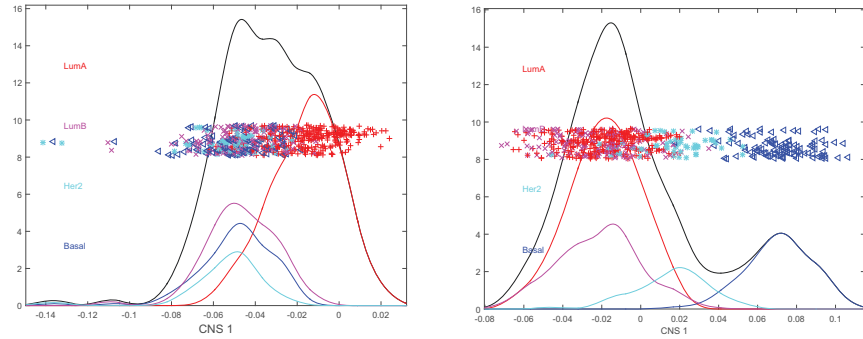


Figure 6: Left: Kernel density estimates of the CNS among GE, CN, RPPA and mutation. The clear separation between Luminal A versus the other subtypes indicates that these four data blocks share a very strong Luminal A property captured in this joint variation component; Right: The CNS from applying AJIVE to the individual matrices of GE, CN, and RPPA. The clear separation indicates that these contain a joint variation component that is consistent with the subtype difference between Basal versus the others.

The clear separation among density estimates suggest that this joint variation component is strongly connected with the subtype difference between Luminal A versus the other subtypes. To quantify this subtype difference, a test is performed using the CNS of this joint component evaluated by the DiProPerm hypothesis test (Wei et al., 2015) based on 100 permutations. Strength of the evidence is usually measured by permutation p-values. However, in this context empirical p-values are frequently zero. Thus a more interpretable measure of strength of the evidence is the DiProPerm z-score. This is 29.32 for this CNS. An area under the receiver operating characteristic (ROC) curve (AUC) (Hanley and McNeil, 1982) of 0.915, is also obtained to reflect the classification accuracy. These numbers confirm the strong Luminal A property shared by these four data types.

A further understanding can be obtained by identifying the feature set of each data type which jointly works with the others in characterizing the Luminal A property. By studying the loading coefficients, important mutation features TP53, TTN and PIK3CA are identified which are well known features from previous studies. Similarly the dominants of GATA3 in RPPA is well known, and is connected with the large GATA3 mutation loading. A less well known result of this analysis is the genes appearing with large GE loadings. Many of these are not dominant in earlier studies, which had focused on subgroup separation, instead of joint behavior.

As noted in the discussion of Step 2 above, all four data types have only one significant joint component. However, the individual components for all of GE, CN and RPPA seem to have 3-way joint components. This is investigated by performing a second AJIVE analysis. In particular, we apply the second and third step to the 3 individual variation matrices from the initial analysis. Notice that all individual matrices are low rank and thus the first step is not necessary. The AJIVE analysis results in one joint variation component which is displayed in the right panel of Figure 6. This joint variation component clearly shows the differences among Basal, HER2 and Luminal subtypes. In particular, a subtype difference between Basal-like versus the others is quantified using the DiProPerm z-score (29.82) and the AUC (0.998). Considering the fact that the AUC of the classification between Basal-like versus the others using all the original separate GE features is 0.999, this single joint component contains almost all the variation information for separating Basal-like from the others. This hierarchical application of AJIVE reveals an important joint component that is specific to GE, CN and RPPA but not to Mutation.

We repeated the analysis using a higher number of individual components, selecting 11 (GE), 16 (CN), 33 (RPPA) and 29 (Mutation) in step 1. We found two joint components with the first joint component similar to the 4-way joint component and the second joint component similar to the 3-way joint component discussed above. The running time of steps 2 and 3 increased to around 10 minutes.

4.2. Spanish Mortality Data

A quite different data set from the Human Mortality Database is studied here, which consists of both Spanish males and females. For each gender data block, there is a matrix of *mortality*, defined as the number of people who died divided by the total, for a given age group and year. Because mortality varies by several orders of magnitude, the \log_{10} of the mortality is studied here. Each row represents an age group from 0 to 95, and each column represents a year between 1908 and 2002. In order to associate the historical events with the variations of mortality, columns (i.e. mortality as a function of age) are considered as the common set of data objects of each gender block. Marron and Alonso (2014) performed analysis on the male block and showed interesting interpretations related to Spanish history. Here we are looking for a deeper analysis which integrates both males and females by exploring joint and individual variation patterns.

AJIVE is applied to the two gender blocks centered by subtracting the mean of each age group. The most interesting AJIVE analysis comes from 3 male and 3 female components. The resulting AJIVE gives 2 joint components and 1 of each individual component. Since the loading matrices provide important information on the effect of

different age groups, BSS analysis together with loading matrices is most informative here.

Figure 7 shows a view of the first joint components for the males (left) and females (right) that is very different from the heat map views used in Section 1.1. While these components are matrices, additional insights come from plotting the rows of the matrices as curves over year (top) and the columns as curves over age (bottom). The curves over year (top) are colored using a heat color scheme, indexing age (black = 0 through red = 40 to yellow = 95 as shown in the vertical color bar on the bottom left). The curves over age (bottom) are colored using a rainbow color scheme (magenta = 1908 through green = 1960 to red = 2002, shown in the horizontal color bar in the top) and use the vertical axis as domain with horizontal axis as range to highlight the fact that these are column vectors. Additional visual cues to the matrix structure are the horizontal rainbow color bar in the top panel, showing that year indexes columns of the data matrix and the vertical heat color bar (bottom) showing that age indexes rows of the component matrix. Because this is a single component, i.e. a rank one approximation of the data, each curve is a multiple of a single eigenvector. The corresponding coefficients are shown on the right. In conventional PCA/SVD terminology, the upper BSS coefficients are called *loadings*, and are in fact the entries of the left eigenvectors (colored using the heat color scale on the bottom). Similarly, the lower coefficients are called *scores* and are the entries of the right eigenvectors, colored using the rainbow bar shown in the top.

The scores plots together with the rows as curves plots in Figure 7 indicate a dramatic improvement in mortality over time for both males and females. The scores plots are bimodal indicating rapid overall improvement in mortality around the 1950s. This is also visible as the steepest part in the rows as curves plot. Thus the first mode of joint variation is driven by overall improvement in mortality. In addition to the overall improvement, the rows as curves and scores plots also show the major mortality events, the global flu pandemic of 1918 and the Spanish Civil war in the late 1930s. The loading plots together with the columns as curves plots present the different impacts of this common variation on different age groups for males and females. The loadings plot for males suggests the improvement in mortality is gradually increasing from older towards younger age groups. In contrast, the female block has a bimodal kernel density estimate of the loadings. This shows that females of child bearing age have received large benefits from improving health care. This effect is similarly visible from comparing the female versus male columns as curves.

The second BSS components of joint variation within each gender are similarly visualized in Figure 8. This common variation reflects the contrast between the years around 1950 and the years around 1980 which can be told from the curves in the left top and the colors in the right bottom subplots in both male and female panels. In the scores plots, the green circles, seen on the left end, represent the years around 1950 when automobile penetration started. And the orange to red circles on the right end correspond to recent years, and much improved car and road safety. The upper left loadings plot of males shows that these automobile events had a stronger influence on the 20-45 males in terms of both larger values and a second peak in the kernel density estimate. Although this contrast can also be seen in the loadings plot of females, it is not as strong as for the male block. The JC2 loadings plots show an interesting outlier,

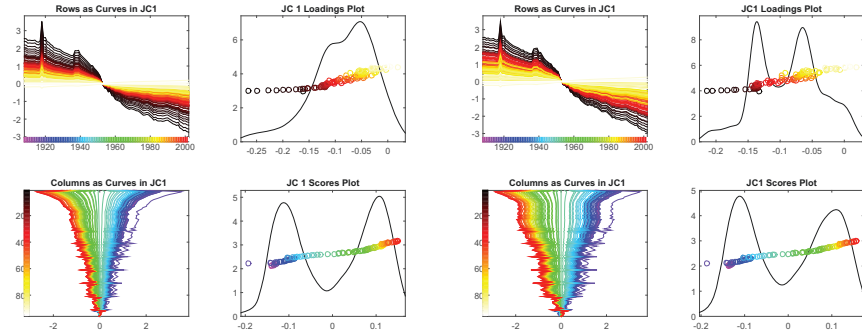


Figure 7: The first BSS joint components of male (left panel) and female (right panel) contain the common modes of variation caused by the overall improvement across different age groups, as can be seen from the scores plots in the right bottom of each panel. The dramatic decrease happened around the 1950s shown in the columns plots. The degree of decrease varies over age groups.

the babies of age zero. We speculate this shows an improvement in post-natal care that coincidentally happened around the same time.

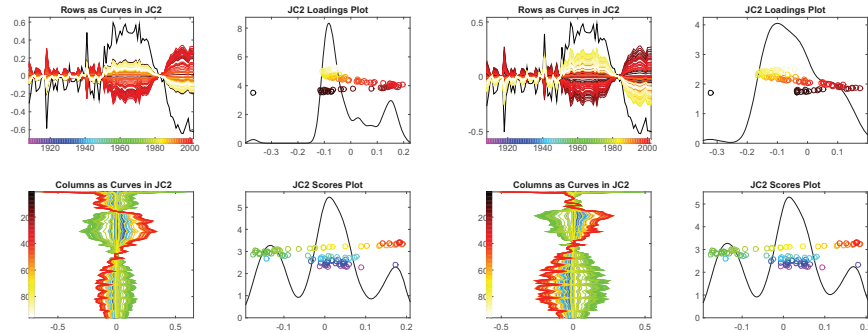


Figure 8: The second joint components of male (left) and female (right) contain the common modes of variation driven by the increase in fatalities caused by automobile penetration and later improvement due to safety improvements. This can be seen from the scores plots in the right bottom. The loadings plots show that this automobile event exerted a significantly stronger impact on the 20-45 males.

Another interesting result comes from the studying first individual components (IC1) of males and females, shown in Figure 9. In the scores plot of males (left), the blue circles stand out from the rest, corresponding to the years of the Spanish civil war when a significant spike can be seen in the rows as curves plot. Young to middle age groups are affected more than the others which can be found in the loadings plot and columns as curves plot. This year variation pattern, however, cannot be detected in the individual variation component of females. The columns as curves plot on the lower left suggest some type of 5-year age rounding effect, which is seen to occur mostly during the earlier years as indicated both in the rows as curves plot and the colors of the peaks in the columns as curves plot. Note that the plot scales show that the

individual female effects are much smaller in magnitude than the male effects.

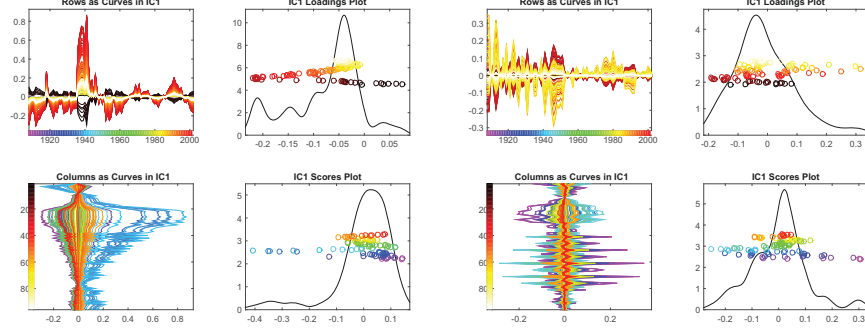


Figure 9: The individual component of male (left) contains the variation driven by the Spanish civil war which can be seen from the blue circles on the right end of the right bottom plot. The Spanish civil war mainly affected the young to middle age males.

5. Optimization perspective

In this section we will investigate how AJIVE compares to PLS, CCA and COBE using the optimization problems that each method is based on. Recall that $X_k, k = 1, \dots, K$ are $(d_k \times n)$ data matrices, with SVD decompositions $X_k = U_{X_k} \Sigma_{X_k} V_{X_k}^T$, where Σ_{X_k} contains no zeros on its diagonal. To be compatible with AJIVE, we will consider these three algorithms in a non-standard configuration using row spaces. In Section 5.1 and 5.2, we assume that the matrices X_k are row centered. We will also use the following notation: for $\vec{a}_1 \in \mathbb{R}^{d_1}, \vec{a}_2 \in \mathbb{R}^{d_2}$

$$\langle \vec{a}_1 X_1, \vec{a}_2 X_2 \rangle = \text{Cov}(\vec{a}_1 X_1, \vec{a}_2 X_2) = \sqrt{\text{Var}(\vec{a}_1 X_1) \text{Var}(\vec{a}_2 X_2)} \text{Corr}(\vec{a}_1 X_1, \vec{a}_2 X_2).$$

5.1. Partial Least Squares

The PLS finds linear combinations of rows of X_1 and X_2 maximizing their sample covariance. More precisely, the PLS identifies a set of pairs of principal vectors, indexed by i , obtained sequentially from the following maximization problems:

$$\begin{aligned} \{\vec{a}_1^{(i)}, \vec{a}_2^{(i)}\} &= \underset{\vec{a}_1 \in \mathbb{R}^{d_1}, \vec{a}_2 \in \mathbb{R}^{d_2}}{\text{argmax}} \quad \langle \vec{a}_1 X_1, \vec{a}_2 X_2 \rangle \\ \text{subject to the constraints: } &\|\vec{a}_1\| = 1, \|\vec{a}_2\| = 1, \\ &\langle \vec{a}_1 X_1, \vec{a}_1^{(j)} X_1 \rangle = 0, \langle \vec{a}_2 X_2, \vec{a}_2^{(j)} X_2 \rangle = 0, j = 1, \dots, i-1. \end{aligned} \quad (10)$$

Unlike AJIVE, the directions from PLS are influenced by both variance within data blocks and correlation between the data blocks. In particular, if the signal strength of the individual structure is sufficiently large it might be mistakenly classified as a joint structure by being found ahead of the real joint structure. This phenomenon can be seen in the analysis of the toy example of Section 1.1 in Appendix B.

5.2. Canonical Correlation Analysis/ Principal Angle Analysis

Similar to PLS, the CCA finds linear combinations of rows of X_1 and X_2 maximizing their sample correlation. In particular, CCA identifies a set of pairs of canonical vectors obtained sequentially from the optimization problem:

$$\begin{aligned} \{\vec{a}_1^{(i)}, \vec{a}_2^{(i)}\} &= \operatorname{argmax}_{\vec{a}_1 \in \mathbb{R}^{d_1}, \vec{a}_2 \in \mathbb{R}^{d_2}} \langle \vec{a}_1 X_1, \vec{a}_2 X_2 \rangle \\ \text{subject to the constraints: } &\|\vec{a}_1 X_1\| = 1, \|\vec{a}_2 X_2\| = 1 \\ &\langle \vec{a}_1 X_1, \vec{a}_1^{(j)} X_1 \rangle = 0, \langle \vec{a}_2 X_2, \vec{a}_2^{(j)} X_2 \rangle = 0, j = 1, \dots, i-1. \end{aligned} \quad (11)$$

This form makes the relationship between (10) and (11) clear and is equivalent to the usual formulation of optimizing the correlation.

There is an important relationship between CCA and PAA (Björck and Golub, 1973), i.e., if $\rho_i = \langle \vec{a}_1^{(i)} X_1, \vec{a}_2^{(i)} X_2 \rangle$ is the i th canonical correlation, $\rho_i = \cos(\theta_i)$, where θ_i is the i th principal angle between row spaces of X_1 and X_2 . The principal vector pairs $\{\vec{x}_{1,i}, \vec{x}_{2,i}\} = \{\vec{a}_1^{(i)} X_1, \vec{a}_2^{(i)} X_2\}$ are often obtained through SVD of $V_{X_1}^T V_{X_2}$. In particular, let $\vec{u}_{X_1,i}, \vec{u}_{X_2,i}$ be the i th left and right singular vectors of $V_{X_1}^T V_{X_2}$. Then, the i th pair of principal vectors are

$$\vec{x}_{1,i} = \vec{u}_{X_1,i}^T V_{X_1}^T, \quad \vec{x}_{2,i} = \vec{u}_{X_2,i}^T V_{X_2}^T.$$

625 An issue with CCA of high-dimensional data is related to the fact that CCA is interested in the canonical vectors \vec{a}_i rather than the principal vectors \vec{x}_i . In particular, when $d_1 > n, d_2 > n$, the values of \vec{a}_i in (11) are not identifiable due to the singularity of $X_1 X_1^T$ and $X_2 X_2^T$. Several approaches have been taken to solve this problem. One approach is to use the Moore-Penrose pseudo inverse to replace the inverse of $X_1 X_1^T$ and $X_2 X_2^T$. A second approach is to add a ridge penalty on $X_1 X_1^T$ and $X_2 X_2^T$ (Vinod, 1976). A third approach called penalized CCA is to add penalty functions on $\{\vec{a}_1^{(i)}, \vec{a}_2^{(i)}\}$, such as an ℓ_1 penalty (Parkhomenko et al., 2007; Lê Cao et al., 2009), an elastic net (Waaijenborg et al., 2008) or a fused lasso (Witten et al., 2009). Another approach called diagonal penalized CCA is to replace $X_1 X_1^T$ and $X_2 X_2^T$ by $\text{diag}(X_1 X_1^T)$ and $\text{diag}(X_2 X_2^T)$ (Parkhomenko et al., 2009; Witten et al., 2009).

635 Another important issue with CCA, which is directly related to AJIVE, is that when $d_1 > n, d_2 > n$, CCA is generally driven by noise. Lee (2007); Samarov (2009); Lee (2016) study the asymptotic behavior of CCA in the high-dimension low sample size context and point out the inconsistency phenomenon in this case. One can solve this issue, like AJIVE and COBE, by replacing X_k by its low rank approximation \tilde{A}_k , $k = 1, 2$. Recall notation from (3). The i th principal vectors are $\vec{p}_i = \tilde{V}_1 \vec{u}_{1,i}$, $\vec{q}_i = \tilde{V}_2 \vec{u}_{2,i}$, where $\vec{u}_{j,i}$ is the i th singular vector of \tilde{U}_i of the SVD of $\tilde{V}_1^T \tilde{V}_2$ respectively.

As discussed in Section 2, AJIVE uses an equivalent principal angle calculation based on SVD of $M = [\tilde{V}_1, \tilde{V}_2]^T = U_M \Sigma_M V_M^T$ (Miao and Ben-Israel, 1992). AJIVE uses the transpose of the i th right singular vector, $V_{M,i}^T$, as the estimated i th basis vector of the joint space, provided that the i th principal angle is smaller than the threshold derived in Section 2.4. Moreover, if the i th principal angle has a value distinct from

other principal angles, then the i th left singular vector of M can be written as $U_{M,i} = [\tilde{u}_{1,i}^T, \tilde{u}_{2,i}^T]^T / \sqrt{2}$. Consequently

$$V_{M,i}^T = \frac{1}{\sigma_{M,i}} U_{M,i}^T M = \frac{1}{\sqrt{2}\sigma_{M,i}} (\tilde{u}_{1,i}^T \tilde{V}_1^T + \tilde{u}_{2,i}^T \tilde{V}_2^T) = \frac{1}{\sqrt{2}\sigma_{M,i}} (\tilde{p}_i^T + \tilde{q}_i^T).$$

This shows that the AJIVE direction $V_{M,i}^T$ is the scaled sum of the i th pair of principal vectors.

645 CCA applied to the low rank approximations \tilde{A}_k and AJIVE are therefore closely related. However, AJIVE provides one joint vector per two distinct principal vectors that by the virtue of being an average should be a better estimate of the joint space than either of the principal vectors. More importantly, AJIVE uses a theoretically sound threshold of the principal angles that allows us to segment individual and joint varia-
650 tion.

The AJIVE formulation allows for a natural extension to multi-block situations. Several approaches of Multiset Canonical Correlation Analysis (mCCA) have been developed as extensions of CCA (Horst, 1961; Kettenring, 1971; Nielsen, 2002). There is no general consensus on which of these extensions is preferable. We point out that
655 AJIVE is closely related to one of the mCCA discussed in Nielsen (2002).

This version of mCCA is defined using the optimization problem for the i th set of canonical vectors $\{\tilde{a}_1^{(i)}, \dots, \tilde{a}_K^{(i)}\}$ and corresponding principal vectors (also called canonical variables) $\{\tilde{a}_1^{(i)} X_1, \dots, \tilde{a}_K^{(i)} X_K\}$:

$$\begin{aligned} \{\tilde{a}_1^{(i)}, \dots, \tilde{a}_K^{(i)}\} &= \operatorname{argmax}_{\tilde{a}_1, \dots, \tilde{a}_K} \sum_{1 \leq k, l \leq K} \langle \tilde{a}_k X_k, \tilde{a}_l X_l \rangle \\ \text{subject to the constraints: } &\sum_{k=1}^K \|\tilde{a}_k X_k\|_2^2 = 1, \\ &\langle \tilde{a}_k X_k, \tilde{a}_k^{(j)} X_k \rangle = 0, \quad k = 1, \dots, K, \quad j = 1, \dots, i-1. \end{aligned} \quad (12)$$

Notice that the constraint in (12) is different than the perhaps more natural $\|\tilde{a}_k X_k\|_2^2 = 1$ for all k .

If the i th singular value corresponding to the AJIVE direction $V_{M,i}^T$ has a value distinct from other singular values in the AJIVE SVD, then calculations similar to the two block case show that the i th basis vector of the joint space from AJIVE

$$V_{M,i}^T = \frac{1}{\sigma_{M,i}} \sum_{k=1}^K \tilde{a}_k^{(i)} X_k$$

is the scaled sum of the corresponding canonical variables. In fact, $V_{M,i}^T$ is the i th flag mean of the row spaces of X_1, \dots, X_K , as defined by Draper et al. (2014), which thus
660 is a building block of AJIVE.

5.3. Common Orthogonal Basis Extraction

Zhou et al. (2016) proposed a compelling optimization problem for finding the

common orthogonal basis (COBE). It is based on iteratively solving

$$\bar{a}_i = \underset{\bar{a}, z_{i,k}, k=1, \dots, K}{\operatorname{argmin}} \sum_{k=1}^K \|\tilde{V}_k z_{i,k} - \bar{a}\|^2 \quad (13)$$

subject to the constraints: $\|\bar{a}\|_2 = 1, \langle \bar{a}, \bar{a}_j \rangle = 0, j = 1, \dots, i-1$.

To compare COBE to AJIVE we first simplify the objective function of (13) to

$$\begin{aligned} \sum_{k=1}^K \|\tilde{V}_k z_{i,k} - \bar{a}_i\|_2^2 &= \sum_{k=1}^K \|\tilde{V}_k z_{i,k}\|_2^2 + K \|\bar{a}_i\|_2^2 - 2 \sum_{k=1}^K \langle \tilde{V}_k z_{i,k}, \bar{a}_i \rangle \\ &= \|z_i\|_2^2 + K \|\bar{a}_i\|_2^2 - 2 z_i^T M \bar{a}_i. \end{aligned}$$

where $z_i = [z_{i,1}, \dots, z_{i,K}]$. If we fix the value of $\|z_i\|$ we see that the solution to the optimization problem (13) is the same as SVD of M with $\bar{a}_i = V_{M,i}$. Moreover this solution is invariant in $\|z_i\|$.

665 Thus the optimization problem (13) gives the same result as AJIVE. However, because AJIVE uses well optimized SVD rather than a heuristic iteration algorithm, AJIVE is much faster than the COBE algorithm. Moreover, COBE lacks any principally based standard on how to choose the threshold for selecting the joint space.

670 To understand why this is a serious issue consider the results of applying COBE to the TCGA data discussed in detail in the next section. To make comparisons fair we provided COBE the same selected first stage ranks for each data block as AJIVE. COBE's default threshold for separating joint and individual structure of 0.01 is too low to find any joint component, even after we raised the selected ranks in the first step to 11 (GE), 16 (CN), 33 (RPPA) and 29 (Mutation).

675 Therefore we tried raising the default threshold 0.01 to 1, in which case COBE fails to finish on our computer due to its inefficient handling of high dimensional data. To see if COBE would run with a smaller three-block data set we removed the highest dimensional data block (CN) from the analysis. When using input ranks 11 (GE), 33 (RPPA), and 29 (Mutation) COBE finished in 6.6 hours returning the joint components of rank 11. This is unreasonable as it selected the largest possible joint space. As a comparison, AJIVE applied to the three-block data and given the same first stage ranks finished in around 6 minutes returning 2 joint components.

Acknowledgements

685 This research was supported in part by the National Science Foundation under Grant No. 1016441, 1512945 and 1633074.

Appendix A. Proofs

Proof of Lemma 2.1. Define the row subspaces respectively for each matrix A_k as $\operatorname{row}(A_k) \subseteq \mathbb{R}^n$. For non-trivial cases, define a subspace $\operatorname{row}(J) \neq \{\vec{0}\}$ as the intersection of the row spaces $\{\operatorname{row}(A_k), k = 1, \dots, K\}$ i.e.

$$\operatorname{row}(J) \triangleq \bigcap_{k=1}^K \operatorname{row}(A_k).$$

For each matrix A_k , two matrices J_k, I_k can be obtained by projection of A_k on $\text{row}(J)$ and its orthogonal complement in the row space $\text{row}(A_k)$. Thus the two matrices satisfy $J_k + I_k = A_k$ and their row subspaces are orthogonal with each other, i.e. $\text{row}(J) \perp \text{row}(I_k)$, $k = 1, \dots, K$. Then the intersection of the row subspaces $\{\text{row}(I_k), k = 1, \dots, K\}$, $\bigcap_{k=1}^K \text{row}(I_k)$, has a zero projection matrix. Therefore, we

have $\bigcap_{k=1}^K \text{row}(I_k) = \{\vec{0}\}$ and have obtained a set of matrices simultaneously satisfying the stated constraints.

On the other hand, it follows from the assumptions that the row space $\text{row}(A_k)$ is spanned by the union of basis vectors of $\text{row}(J_k)$ and $\text{row}(I_k)$, which indicates

$$\text{row}(J) = \bigcap_{k=1}^K \text{row}(A_k).$$

Accordingly, the matrices J_k and I_k for $k = 1, \dots, K$ are also uniquely defined. \square

Proof of Lemma 2.3. Let P_1 and P_2 be the projection matrices onto the individually perturbed joint row spaces. And let P be the projection matrix onto the common joint row space J . Thus, we have

$$\sin \theta = \|(I - P_1)P_2\| \quad (\text{A.1})$$

$$\leq \|(I - P_1)(I - P)P_2\| + \|(I - P_1)PP_2\| \quad (\text{A.2})$$

$$\leq \|(I - P_1)(I - P)\| \|(I - P)P_2\| + \|(I - P_1)P\| \|PP_2\| \quad (\text{A.3})$$

in which $\|(I - P_1)P\| = \sin \theta_{1,1}$, $\|(I - P_1)(I - P)\| = \cos \theta_{1,1}$, $\|(I - P_2)P\| = \sin \theta_{2,1}$ and $\|(I - P_2)(I - P)\| = \cos \theta_{2,1}$. Therefore,

$$\sin \phi \leq \cos \theta_{1,1} \sin \theta_{2,1} + \sin \theta_{1,1} \cos \theta_{2,1} = \sin(\theta_{1,1} + \theta_{2,1}).$$

\square

Proof of Lemma 2.4. Notation from (5) and (7) is used here. For each singular value $\sigma_{M,i}$, it can be formulated as a sequential optimization problem i.e

$$\sigma_{M,i}^2 = \max_Q \|MQ\|_F^2 = \max_Q \sum_{k=1}^K \|\tilde{V}_1^T Q\|_F^2,$$

where Q is a rank 1 projection matrix that is orthogonal to the previous $i - 1$ optima i.e. Q_1, \dots, Q_{i-1} . The Q that maximizes the Frobenius norm of MQ is denoted as Q_i .

For an arbitrary component in the theoretical joint score subspace $\text{row}(J)$, write its projection matrix as $P_J^{(1)}$. The Frobenius norm of M projected onto $P_J^{(1)}$ is

$$\|MP_J^{(1)}\|_F^2 = \left\| \begin{bmatrix} \tilde{V}_1^T P_J^{(1)} \\ \vdots \\ \tilde{V}_K^T P_J^{(1)} \end{bmatrix} \right\|_F^2 \geq \left\| \begin{bmatrix} \cos \theta_1 \\ \vdots \\ \cos \theta_K \end{bmatrix} \right\|_F^2 = \sum_{k=1}^K \cos^2 \theta_k \quad (\text{A.4})$$

Considering the mechanism of SVD, $\sigma_{M,1}^2$ is the maximal norm obtained from the optimal projection matrix $Q_1 \subseteq \bigcup_{k=1}^K \text{row}(\tilde{A}_k) \subseteq \mathbb{R}^n$. If all \tilde{A}_k contain all components obtained by noise perturbation of the common row space $\text{row}(J)$, then we have

$$\sigma_{M,1}^2 \geq \|MP_J^{(1)}\|_F^2 \geq \sum_{k=1}^K \cos^2 \theta_k$$

700 to be considered as a component of the joint score subspace.

This argument can be applied sequentially. For the $Q_2 \in Q_1^\perp \cap \{\bigcup_{k=1}^K \text{row}(\tilde{A}_k)\}$, there exist a non-empty joint subspace ($\subseteq \text{row}(J)$) such that all $Q_1^\perp \cap \text{row}(\tilde{A}_k)$ contain perturbed directions of a joint component other than the one above. Therefore this joint component with projection matrix $P_J^{(2)}$ should have

$$\sigma_{M,2}^2 \geq \|MP_J^{(2)}\|_F^2 \geq \sum_{k=1}^K \cos^2 \theta_k.$$

Thus the singular values corresponding to the joint components satisfies (8) and this procedure can continue through at least r_J steps. \square

Appendix B. Details of the toy example

Section 1.1 introduces a toy example of two data blocks, X (100×100) and Y (10000 \times 100), with patterns corresponding to joint and individual structures. For details see Figure 2.

A naive attempt at integrative analysis can be done by concatenating X and Y on columns and performing a singular value decomposition on this concatenated matrix. Figure B.10 shows the results for 3 choices of rank. The rank 2 approximation essentially captures the joint variation component and the individual variation component of X , but the Y components are hard to interpret. The bottom 2000 rows show the joint variation but the top half of Y reveals signal from the individual component of X . One might hope that the Y individual components would show up in the rank 3 and rank 4 approximations. However, because the noise in the X matrix is so large, a random noise component from X dominates the Y signal, so the important latter component disappears from this low rank representation unlike the AJIVE result in Figure 3. In this example, this naive approach completely fails to give a meaningful joint analysis.

Figure B.11 presents the PLS approximations with different numbers of components selected. PLS completely fails to separate joint and individual components. Instead it provides mixtures of the joint, and some of the individual components. Increasing the rank of the PLS approximation only includes more noise.

The Lock et al. (2013) method, called JIVE here, is applied to this toy data set. The left panel of Figure B.12 shows a reasonable approximation of the total signal variation within each data block. However, the Lock et al. (2013) method gives rank 2 approximations to the joint matrices shown in the middle panel. The approximation consists of the real joint component together with the individual component of X . Following this, the approximation of the X individual matrix is a zero matrix and a

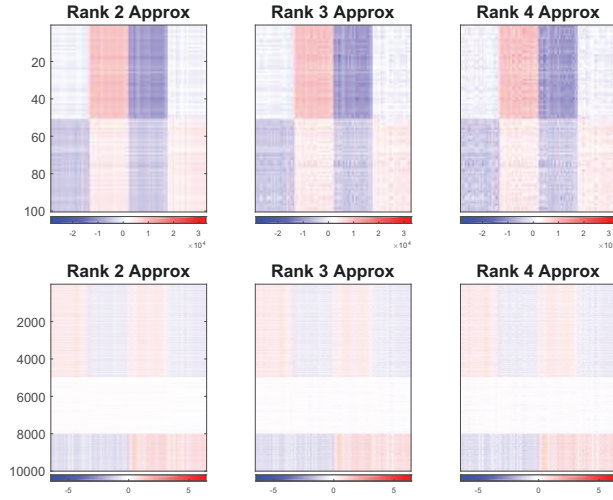


Figure B.10: Shows the concatenation SVD approximation of each block for rank 2 (left), 3 (center) and 4 (right). Although block X has a relatively accurate approximation when the rank is chosen as 2, the individual pattern in block Y has never been captured due to the heterogeneity between X and Y .

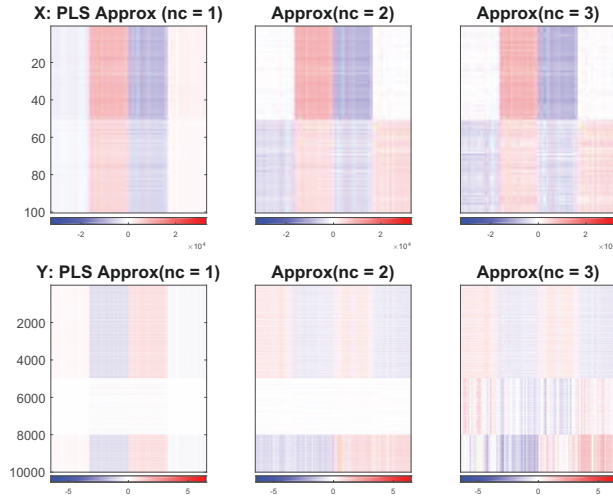


Figure B.11: PLS approximations of each block for numbers of components as 1 (left), 2 (center) and 3 (right). PLS fails to distinguish the joint and individual variation structure.

730 wrong approximation of the Y individual matrix is obtained shown in the top half of the right panel. We speculate that failure to correctly apportion the joint and individual variation is caused by the fact that the individual spaces are correlated. Because of relying on a permutation test JIVE's algorithm does not handle correlated individual signals very well.

735 We finally remark that the Zhou et al. (2016) method COBE correctly segments the toy example. However it takes significantly (42 times) longer time than AJIVE to do SO.

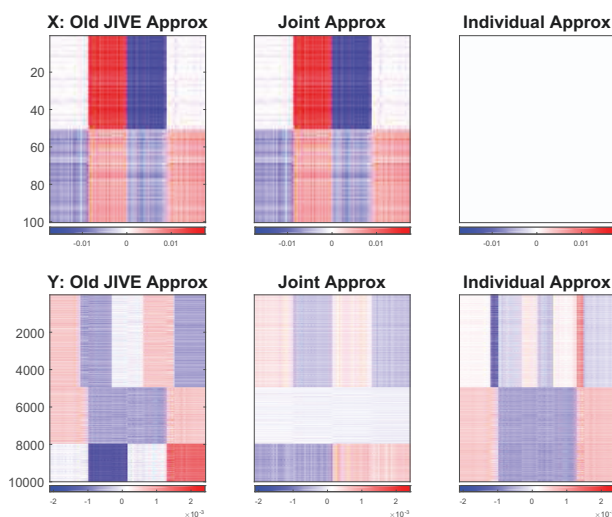


Figure B.12: The Lock et al. (2013) JIVE method approximation of the data blocks X and Y in the toy example are shown in the first panel of figures. The joint matrix approximations (middle panel) incorrectly contain the individual component of X caused by the problematic algorithm and inappropriate normalization.

References

- Abdi, H., Williams, L.J., Valentin, D., 2013. Multiple factor analysis: principal component analysis for multitable and multiblock data sets. *Wiley Interdisciplinary reviews: computational statistics* 5, 149–179.
- 740 Björck, Å., Golub, G.H., 1973. Numerical methods for computing angles between linear subspaces. *Mathematics of computation* 27, 579–594.
- Cai, T.T., Zhang, A., 2016. Rate-optimal perturbation bounds for singular subspaces with applications to high-dimensional statistics. *arXiv preprint arXiv:1605.00353*.
- 745 Ciriello, G., Gatza, M.L., Beck, A.H., Wilkerson, M.D., Rhie, S.K., Pastore, A., Zhang, H., McLellan, M., Yau, C., Kandath, C., et al., 2015. Comprehensive molecular portraits of invasive lobular breast cancer. *Cell* 163, 506–519.

- Draper, B., Kirby, M., Marks, J., Marrinan, T., Peterson, C., 2014. A flag representation for finite collections of subspaces of mixed dimensions. *Linear Algebra and its Applications* 451, 15–32.
- 750 Hanafi, M., Kohler, A., Qannari, E.M., 2011. Connections between multiple co-inertia analysis and consensus principal component analysis. *Chemometrics and intelligent laboratory systems* 106, 37–40.
- Hanley, J.A., McNeil, B.J., 1982. The meaning and use of the area under a receiver operating characteristic (roc) curve. *Radiology* 143, 29–36.
- 755 Horst, P., 1961. Relations among m sets of measures. *Psychometrika* 26, 129–149.
- Hotelling, H., 1936. Relations between two sets of variates. *Biometrika* 28, 321–377.
- Jere, S., Dauwels, J., Asif, M.T., Vie, N.M., Cichocki, A., Jaillet, P., 2014. Extracting commuting patterns in railway networks through matrix decompositions, in: *Control Automation Robotics & Vision (ICARCV)*, 2014 13th International Conference on, IEEE. pp. 541–546.
- 760 Jordan, C., 1875. Essai sur la géométrie à n dimensions. *Bulletin de la Société mathématique de France* 3, 103–174.
- Kettenring, J.R., 1971. Canonical analysis of several sets of variables. *Biometrika* , 433–451.
- 765 Kotz, S., Nadarajah, S., 2004. *Multivariate t -distributions and their applications*. Cambridge University Press.
- Kühnle, O., 2011. Integration of multiple high-throughput data-types in cancer research. Ph.D. thesis. Ludwig Maximilian University of Munich.
- Kuligowski, J., Pérez-Guaita, D., Sánchez-Illana, Á., León-González, Z., de la Guardia, M., Vento, M., Lock, E.F., Quintás, G., 2015. Analysis of multi-source metabolomic data using joint and individual variation explained (jive). *Analyst* .
- 770 Lê Cao, K.A., Martin, P.G., Robert-Granié, C., Besse, P., 2009. Sparse canonical methods for biological data integration: application to a cross-platform study. *BMC bioinformatics* 10, 34.
- 775 Lee, M.H., 2007. Continuum direction vectors in high dimensional low sample size data. Ph.D. thesis. University of North Carolina at Chapel Hill.
- Lee, S., 2016. High-dimension, low sample size asymptotics of canonical correlation analysis. *arXiv preprint arXiv:1609.02992* .
- Lock, E.F., Dunson, D.B., 2013. Bayesian consensus clustering. *Bioinformatics* 29, 2610–2616.
- 780 Lock, E.F., Hoadley, K.A., Marron, J., Nobel, A.B., 2013. Joint and individual variation explained (JIVE) for integrated analysis of multiple data types. *The Annals of Applied Statistics* 7, 523.

- Löfstedt, T., Hoffman, D., Trygg, J., 2013. Global, local and unique decompositions in onpls for multiblock data analysis. *Analytica chimica acta* 791, 13–24.
- Marron, J.S., Alonso, A.M., 2014. Overview of object oriented data analysis. *Biometrical Journal* 56, 732–753.
- Miao, J., Ben-Israel, A., 1992. On principal angles between subspaces in R^n . *Linear algebra and its applications* 171, 81–98.
- Mo, Q., Wang, S., Seshan, V.E., Olshen, A.B., Schultz, N., Sander, C., Powers, R.S., Ladanyi, M., Shen, R., 2013. Pattern discovery and cancer gene identification in integrated cancer genomic data. *Proceedings of the National Academy of Sciences* 110, 4245–4250.
- Network, C.G.A., et al., 2012. Comprehensive molecular portraits of human breast tumours. *Nature* 490, 61–70.
- Nielsen, A.A., 2002. Multiset canonical correlations analysis and multispectral, truly multitemporal remote sensing data. *IEEE transactions on image processing* 11, 293–305.
- O’Connell, M.J., Lock, E.F., 2016. R. JIVE for exploration of multi-source molecular data. *Bioinformatics* 32, 2877 – 2879.
- Parkhomenko, E., Tritchler, D., Beyene, J., 2007. Genome-wide sparse canonical correlation of gene expression with genotypes, in: *BMC proceedings*, BioMed Central. p. S119.
- Parkhomenko, E., Tritchler, D., Beyene, J., et al., 2009. Sparse canonical correlation analysis with application to genomic data integration. *Statistical Applications in Genetics and Molecular Biology* 8, 1–34.
- Ray, P., Zheng, L., Lucas, J., Carin, L., 2014. Bayesian joint analysis of heterogeneous genomics data. *Bioinformatics* 30, 1370–1376.
- Samarov, D.V., 2009. The analysis and advanced extensions of canonical correlation analysis. Ph.D. thesis. University of North Carolina at Chapel Hill.
- Schouteden, M., Van Deun, K., Pattyn, S., Van Mechelen, I., 2013. Sca with rotation to distinguish common and distinctive information in linked data. *Behavior research methods* 45, 822–833.
- Schouteden, M., Van Deun, K., Wilderjans, T.F., Van Mechelen, I., 2014. Performing DISCO-SCA to search for distinctive and common information in linked data. *Behavior Research Methods* 46, 576–587.
- Smilde, A.K., Westerhuis, J.A., de Jong, S., 2003. A framework for sequential multi-block component methods. *Journal of chemometrics* 17, 323–337.
- Stewart, G., Sun, J.g., 1990. *Matrix Perturbation Theory*. Computer science and scientific computing, Academic Press.

- Trygg, J., Wold, S., 2003. O2-PLS, a two-block (X–Y) latent variable regression (LVR) method with an integral OSC filter. *Journal of chemometrics* 17, 53–64.
- Vinod, H.D., 1976. Canonical ridge and econometrics of joint production. *Journal of econometrics* 4, 147–166.
- 825 Waaijenborg, S., de Witt Hamer, P.V., Zwinderman, A.H., et al., 2008. Quantifying the association between gene expressions and dna-markers by penalized canonical correlation analysis. *Statistical Applications in Genetics and Molecular Biology* 7, 3.
- 830 Wedin, P.Å., 1972. Perturbation bounds in connection with singular value decomposition. *BIT Numerical Mathematics* 12, 99–111.
- Wei, S., Lee, C., Wichers, L., Marron, J., 2015. Direction–projection–permutation for high dimensional hypothesis tests. *Journal of Computational and Graphical Statistics* 0, 0–0.
- 835 Westerhuis, J.A., Kourti, T., MacGregor, J.F., 1998. Analysis of multiblock and hierarchical pca and pls models. *Journal of chemometrics* 12, 301–321.
- Witten, D.M., Tibshirani, R., Hastie, T., 2009. A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics* 10, 515 – 534.
- Wold, H., 1985. Partial least squares. *Encyclopedia of statistical sciences* .
- 840 Wold, S., Geladi, P., Esbensen, K., Öhman, J., 1987. Multi-way principal components- and pls-analysis. *Journal of chemometrics* 1, 41–56.
- Wold, S., Kettaneh, N., Tjessem, K., 1996. Hierarchical multiblock pls and pc models for easier model interpretation and as an alternative to variable selection. *Journal of chemometrics* 10, 463–482.
- 845 Yang, Z., Michailidis, G., 2015. A non-negative matrix factorization method for detecting modules in heterogeneous omics multi-modal data. *Bioinformatics* 32, 1 – 8.
- Yu, Q., Risk, B.B., Zhang, K., Marron, J., 2017. Jive integration of imaging and behavioral data. *NeuroImage* 152, 38–49.
- 850 Zhang, Y., Zhou, G., Jin, J., Wang, X., Cichocki, A., 2015. Ssvep recognition using common feature analysis in brain–computer interface. *Journal of neuroscience methods* 244, 8–15.
- 855 Zhou, G., Cichocki, A., Zhang, Y., Mandic, D., 2016. Group component analysis for multiblock data: Common and individual feature extraction. *IEEE Transactions on Neural Networks and Learning Systems* 17, 2426 – 2439.