#### **Lecture 3** Dimension Reduction for Multi-Omics Data

April 18, 2023



Instructor: Jack Pattee

# Outline

#### **Dimension Reduction**

• Principal components analysis

#### Multi-omics Dimension Reduction

- Stacked PCA
- Joint and Individual Variation Explained
- Method Comparison
- Clustering
- Association
- Prediction

## **Dimension Reduction**

- Dimension reduction seeks to explain the variation in a set of 'independent' or predictor variables using a set of derived features.
  - Often in the case of high-dimensional data.
- Considered a form of unsupervised machine learning; i.e., it does not consider an outcome variable.
- Many uses for dimension reduction
  - Clustering: grouping subjects or variables in the dataset.
  - Feature engineering: generating a set of derived features for subsequent use, i.e., in a predictive modeling application.
  - Visualization of high dimensional data.
  - Etc.

#### Principal Components Analysis

- A popular method for dimension reduction is principal components analysis (PCA).
- PCA projects the dataset into a new vector space with an orthonormal basis.
- If we consider our features to be encoded in the *n*-sample by *p*-feature matrix **X**, we can think of the principal components representation as closely related to the singular value decomposition of **X**.

## Singular Value Decomposition

- According to SVD, any  $n \times p$  matrix **X** can be decomposed into:  $X = U\Sigma V^T$ 
  - U is an  $n \times n$  matrix where columns are orthonormal.
  - $\Sigma$  is an  $n \times p$  rectangular diagonal matrix with entries equal to the square root of the eigenvalues of  $X^T X$ .
  - V is a  $p \times p$  matrix with orthonormal columns.
- Note:  $X^T X = \cdots = V^T \Sigma^2 V$ , which satisfies the conditions for an eigendecomposition.
  - Thus, columns of V are equivalent to eigenvectors of the covariance matrix of X
- Define  $W^T = \Sigma V^T$
- This gives us the PCA decomposition: T = XW.

## PCA

#### Consider representation T = XW.

- T is the principal component representation of our data, i.e., an  $n \times p$  matrix where our original data has been projected into an orthonormal basis space. Each column is a principal component.
- T has two special properties:
  - Each column of **T** is orthogonal to all other columns.
  - Each column of **T** explains more variation than all subsequent columns and less variation than all preceding columns.
- W is the *loading matrix* of the principal component representation.  $PC1 = w_{11}X_1 + w_{12}X_2 + \dots + w_{1p}X_p$   $PC2 = w_{21}X_1 + w_{22}X_2 + \dots + w_{2p}X_p$   $PC3 = w_{31}X_1 + w_{32}X_2 + \dots + w_{3p}X_p$

#### PCA

- If we want to generate a 1-dimensional representation of our data X in such a way that preserves the most variation in our data, this is done by projecting the data onto the first eigenvector of the covariance matrix.
  - i.e., *XW*<sub>1</sub> has the largest sample variance among all (normalized) linear combinations of the columns of X.
- We can reduce the number of eigenvectors we use to reconstruct our data, thus projecting the data onto only those top few eigenvectors of the covariance matrix:

$$T_L = XW_L$$

Where  $W_L$  retains only the first L eigenvectors.



#### **Proportion of total variance explained by each PC**

**Cumulative proportion of variance explained** 

Principal component

Principal component



Elements of Statistical Learning (2<sup>nd</sup> Ed.), Hastie, Tibshirani, and Friedman 2009, Chapter 3



Elements of Statistical Learning (2<sup>nd</sup> Ed.), Hastie, Tibshirani, and Friedman 2009, Chapter 14

## Dimension Reduction for Multi-omics

- According to Cantini et al: "Multiomics integrative approaches should be able to capture not only signals shared by all omics data, but also those emerging from the complementarity of various omics data".
- Want a way to model variation shared among datatypes and unique to each datatype.
- Informatively combining information from multiple datatypes allows for more precise analysis.

## Multi-omics

- Consider the multi-omics context. We have *n* individuals where we have measured omic values for some set of *k* omics types. So, say,  $X_1, \ldots, X_k$ , for  $k \ge 2$ , where  $X_i$  is  $p_i \times n$ .
  - Variables are rows, subjects are columns.
  - $p_i$  can vary widely according to study and data type; transcriptomic data may have  $p_i \sim 20,000$ , whereas miRNA may have  $p_i \sim 1,000$
- We want a dimension-reduced representation of the multi-omic data.



Cantini et al, Nature Communications 2021

## Multi-omics

- Rows of the factor matrix **F** are called 'factors'.
- Columns the weight matrix are called 'metagenes'.
- Factors are projections onto the sample space, whereas metagenes are projections onto the omics space.
- The factor matrix can be used to cluster samples, whereas columns of weight matrix can be used to extract markers (i.e., selecting the top ranked genes).

#### Multi-omics

• Consider 'stacking' the k matrices:

$$\boldsymbol{X} = \begin{pmatrix} \boldsymbol{X}_1 \\ \vdots \\ \boldsymbol{X}_k \end{pmatrix}$$

- Could simply perform PCA on the stacked matrix X: this is termed 'consensus PCA' (Wold, Kettaneh, and Tjessem, 1996). Consensus PCA is leveraged by popular 'iCluster' method (Shen, Olshen, and Ladanyi, 2009).
- However, such an approach cannot measure the so-called 'joint' and 'individual' structure.
  - Broadly: this is the distinction between whether or not a dimension reduction method estimates omics-specific factors.
- For example: could be certain patterns unique to miRNA and transcriptomic data, and certain patterns shared between the two.

### Joint and Individual Variation Explained

- One example of a method that estimates both individual and joint variation is JIVE (Lock et al, 2013).
- Idea is to decompose X into joint and individual variation, where matrix  $J_i$  encodes the contribution of omics datatype I to the joint structure and  $A_i$  encodes the individual variation.
- Error terms indicate that we use a reduced representation (i.e., limit the number of ranks used in the reconstruction).

$$X_{1} = J_{1} + A_{1} + \varepsilon_{1}$$
$$\vdots$$
$$X_{k} = J_{k} + A_{k} + \varepsilon_{k},$$



Lock et al, Ann Appl Stat 2013

# JIVE

- JIVE can be conceptualized as a form of multidimensional PCA.
- X has an SVD decomposition, and thus a PCA representation X=WU.
- Thus, the rank r joint matrix J can be written as the product of a  $p \times r$  loading matrix W and an  $r \times n$  factor matrix U:

$$\boldsymbol{W} = \begin{pmatrix} \boldsymbol{W}_1 \\ \vdots \\ \boldsymbol{W}_k \end{pmatrix}$$

Where each  $W_i$  gives the loadings of the joint structure corresponding to datatype *i*.

• The rank *r* individual structure matrix  $A_i$  for  $X_i$  can be written as  $V_i U_i$ , where  $V_i$  is a  $p_i \times r_i$  loading matrix and  $U_i$  is an  $r_i \times n$  factor matrix.



• This gives us the full model as follows:

$$X_1 = W_1 U + V_1 U_1 + R_1$$
  
$$\vdots$$
  
$$X_k = W_k U + V_k U_k + R_k$$

# Applications

- Now that we have a multi-omic dimension reduction, what can we do with it?
  - 1. Cluster samples.
  - 2. Predict clinical outcome.
  - 3. Associate factors with clinical variables of interest (age, gender).
- Review paper (Cantini et al) compares nine methods: intNMF, JIVE, MCIA, RGCCA, iCluster, MOFA, tICA, MSFA, and data fusion.

# Single-Cell Clustering

- Cantini et al used factors derived from multi-omic dimension reduction methods to cluster samples in single cell data.
  - These methods were originally designed for bulk data, but characterizing their performance on single-cell data is of interest.
- Two methods, intNMF and iCluster, intrinsically perform clustering; for the other seven methods, k-means consensus clustering was applied to the factor matrix.
- Two 'omics types: scRNA-seq and scATAC-seq, which measure gene expression and chromatin accessibility, respectively.
- Applied to 206 cells from three cancer cell lines.



### Application to TCGA Data

- Analyzed data from The Cancer Genome Atlas on ten different cancer types.
- Omics data types were gene expression, DNA methylation and miRNA expression.

# TCGA Clustering

- Used clinical subtyping to assess clustering performance of jDR methods.
  - Note: these cannot be considered a 'ground truth' assessment.
- For BRCA, compared clustering to two subtypes: ER/PR/HER-2, and COCA classification.
  - COCA: integrative subtyping obtained by separately clustering different 'omics types and performing a consensus clustering on the shared results.
- For methods that do not automatically perform clustering, conducted k-means using the first four derived factors.
- JIVE has best performance based on clinical subtyping.



ASA' Section on Statistics in Genomics and Genetics

# Predicting Survival

- Used first ten factors to predict survival via Cox regression.
- Examined the number of derived features associated with survival.
- Appears that the number of associated factors depends more on the cancer type than the dimension reduction approach.
- RGCCA, MCIA, and JIVE showed most promising results for those cancers where predicting survival was most difficult.





#### Association with Clinical Factors

- Assess the association of dimension-reduced factors with clinical annotations.
  - May be able to determine if one factor is a proxy for some clinical feature.
- Performed analysis in TCGA data, using three omics datatypes.
- Investigated four variables: 'age', 'days to new tumor', 'gender', and 'neo-adjuvant therapy administration'.
- Investigated selectivity of methods, i.e., ability to generate a one-toone map between reduced factors and clinical annotations.
- RGCCA, MCIA, and MOFA are the best-performing algorithms overall.



#### References

- Elements of Statistical Learning (2<sup>nd</sup> Ed.), Hastie, Tibshirani, and Friedman 2009.
- Cantini, Laura, et al. "Benchmarking joint multi-omics dimensionality reduction approaches for the study of cancer." *Nature communications* 12.1 (2021): 124.
- Wold, Svante, Nouna Kettaneh, and Kjell Tjessem. "Hierarchical multiblock PLS and PC models for easier model interpretation and as an alternative to variable selection." *Journal of chemometrics* 10.5-6 (1996): 463-482.
- Shen, Ronglai, et al. "Integrative subtype discovery in glioblastoma using iCluster." PloS one 7.4 (2012): e35236.
- Lock, Eric F., et al. "Joint and individual variation explained (JIVE) for integrated analysis of multiple data types." *The annals of applied statistics* 7.1 (2013): 523.