

INTEGRATED PRINCIPAL COMPONENTS ANALYSIS (IPCA)

Tiffany Tang

University of California, Berkeley

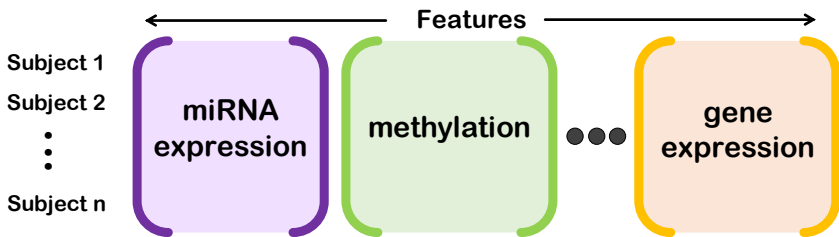
July 30, 2019

Joint Statistical Meetings

Joint work with Genevera Allen (Rice University)

MOTIVATION

- Data Integration + Unsupervised Learning
- Applications
 - Integrative genomics
 - Multi-modal imaging
 - Multi-sensor data
- Want to find the **joint** patterns which are common among all of the datasets



Multiblock PCA

- Concatenated PCA
 - *Westerhuis et al. (1998), Wold et al. (1996)*
- Multiple Factor Analysis
 - *Escofier and Pages (1994), Abdi et al. (2013)*

Matrix Factorization Methods

- Joint and Individual Variation Explained (JIVE)
 - *Lock et al. (2013)*
- Coupled Matrix Factorizations (CMF)
 - *Singh and Gordon (2008), Acar et al. (2014)*

$$\operatorname{argmin}_{\mathbf{U}, \mathbf{V}_1, \dots, \mathbf{V}_K} \sum_{k=1}^K \|\mathbf{X}_k - \mathbf{U}\mathbf{V}_k^T\|_F^2$$

OVERVIEW OF INTEGRATED PCA (IPCA)

Objective: Extend a model-based PCA to integrated data

- Exploratory Data Analysis
- Joint Pattern Recognition
- Visualization

Key tool: matrix-variate normal distribution

Why? PCA can be viewed as maximizing the sample covariance

$$\hat{\Delta} = \frac{1}{n} \mathbf{X}^T \mathbf{X}$$

Advantages:

- A unifying framework for the multiblock PCA family
- **U**'s and **V**'s are orthogonal, ordered, and nested
- Convenient visualizations
- Nice theoretical properties
 - Global solution - important for interpretability and reproducibility
 - Provable guarantees

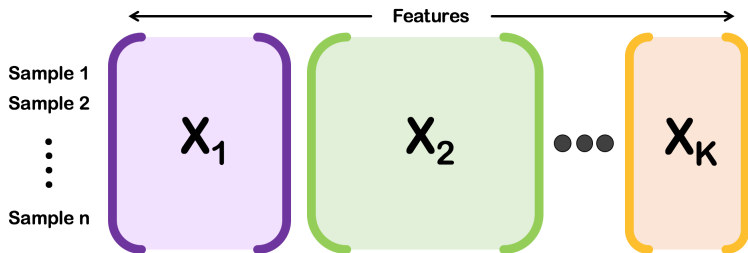
1. iPCA Model
2. Case Study: Alzheimer's Disease

IPCA MODEL

Given datasets $\mathbf{X}_1, \dots, \mathbf{X}_K$, assume that each dataset \mathbf{X}_k arises from the matrix-variate normal model:

$$\mathbf{X}_k \sim N_{n,p_k}(\mathbf{M}_k, \Sigma \otimes \Delta_k).$$

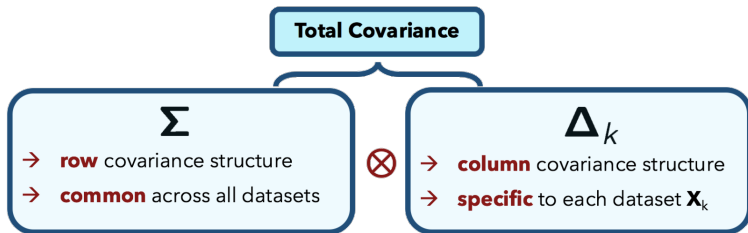
$$\Sigma \otimes \Delta_k = \begin{bmatrix} \sigma_{11} \Delta_k & \cdots & \sigma_{1n} \Delta_k \\ \vdots & \ddots & \vdots \\ \sigma_{n1} \Delta_k & \cdots & \sigma_{nn} \Delta_k \end{bmatrix}$$



IPCA MODEL

Given datasets $\mathbf{X}_1, \dots, \mathbf{X}_K$, assume that each dataset \mathbf{X}_k arises from the matrix-variate normal model:

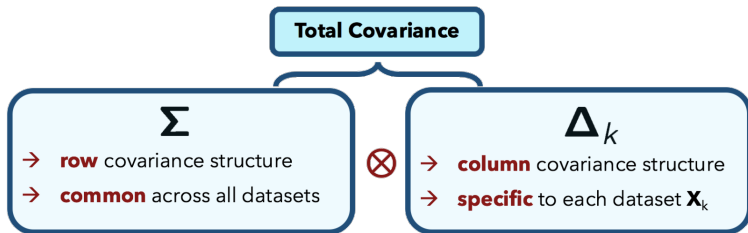
$$\mathbf{X}_k \sim N_{n,p_k}(\mathbf{M}_k, \Sigma \otimes \Delta_k).$$



IPCA MODEL

Given datasets $\mathbf{X}_1, \dots, \mathbf{X}_K$, assume that each dataset \mathbf{X}_k arises from the matrix-variate normal model:

$$\mathbf{X}_k \sim N_{n,p_k}(\mathbf{M}_k, \Sigma \otimes \Delta_k).$$



Equivalently, can rewrite Kronecker product covariance model as

$$\left[\mathbf{X}_1 \Delta_1^{-1/2}, \dots, \mathbf{X}_K \Delta_K^{-1/2} \right]_{.j} \stackrel{iid}{\sim} N(\mathbf{o}, \Sigma)$$

1. iPCA Model: For each $k = 1, \dots, K$, assume

$$\mathbf{X}_k \sim N_{n,p_k}(\mathbf{1}_n \mu_k^T, \Sigma \otimes \Delta_k)$$

- Σ is the **common row** covariance structure
- Δ_k is the **separate column** covariance structure

1. iPCA Model: For each $k = 1, \dots, K$, assume

$$\mathbf{X}_k \sim N_{n,p_k}(\mathbf{1}_n \mu_k^T, \Sigma \otimes \Delta_k)$$

- Σ is the **common row** covariance structure
- Δ_k is the **separate column** covariance structure

2. **Estimate** covariances $\Sigma, \Delta_1, \dots, \Delta_K$ to obtain $\hat{\Sigma}, \hat{\Delta}_1, \dots, \hat{\Delta}_K$

⋮

1. iPCA Model: For each $k = 1, \dots, K$, assume

$$\mathbf{X}_k \sim N_{n,p_k}(\mathbf{1}_n \mu_k^T, \Sigma \otimes \Delta_k)$$

- Σ is the **common row** covariance structure
- Δ_k is the **separate column** covariance structure

2. Estimate covariances $\Sigma, \Delta_1, \dots, \Delta_K$ to obtain $\hat{\Sigma}, \hat{\Delta}_1, \dots, \hat{\Delta}_K$

⋮

3. Maximize covariances

1. iPCA Model: For each $k = 1, \dots, K$, assume

$$\mathbf{X}_k \sim N_{n,p_k}(\mathbf{1}_n \mu_k^T, \Sigma \otimes \Delta_k)$$

- Σ is the **common row** covariance structure
- Δ_k is the **separate column** covariance structure

2. **Estimate** covariances $\Sigma, \Delta_1, \dots, \Delta_K$ to obtain $\hat{\Sigma}, \hat{\Delta}_1, \dots, \hat{\Delta}_K$

⋮

3. **Maximize** covariances

$\mathbf{U} \leftarrow$ eigenvectors of $\hat{\Sigma}$ = joint patterns

$\mathbf{V}_k \leftarrow$ eigenvectors of $\hat{\Delta}_k$ = individual patterns

1. iPCA Model: For each $k = 1, \dots, K$, assume

$$\mathbf{X}_k \sim N_{n,p_k}(\mathbf{1}_n \mu_k^T, \Sigma \otimes \Delta_k)$$

- Σ is the **common row** covariance structure
- Δ_k is the **separate column** covariance structure

2. **Estimate** covariances $\Sigma, \Delta_1, \dots, \Delta_K$ to obtain $\hat{\Sigma}, \hat{\Delta}_1, \dots, \hat{\Delta}_K$

⋮

3. **Maximize** covariances

iPC Scores: $\mathbf{U} \leftarrow$ eigenvectors of $\hat{\Sigma} =$ joint patterns

iPC Loadings: $\mathbf{V}_k \leftarrow$ eigenvectors of $\hat{\Delta}_k =$ individual patterns

1. iPCA Model: For each $k = 1, \dots, K$, assume

$$\mathbf{X}_k \sim N_{n,p_k}(\mathbf{1}_n \mu_k^T, \Sigma \otimes \Delta_k)$$

- Σ is the **common row** covariance structure
- Δ_k is the **separate column** covariance structure

2. **Estimate** covariances $\Sigma, \Delta_1, \dots, \Delta_K$ to obtain $\hat{\Sigma}, \hat{\Delta}_1, \dots, \hat{\Delta}_K$

⋮

3. **Maximize** covariances

iPC Scores: $\mathbf{U} \leftarrow$ eigenvectors of $\hat{\Sigma}$ = joint patterns

iPC Loadings: $\mathbf{V}_k \leftarrow$ eigenvectors of $\hat{\Delta}_k$ = individual patterns

4. **Visualize** dominant joint patterns by plotting iPC scores \mathbf{U}

1. iPCA Model: For each $k = 1, \dots, K$, assume

$$\mathbf{X}_k \sim N_{n,p_k}(\mathbf{1}_n \mu_k^T, \Sigma \otimes \Delta_k)$$

- Σ is the **common row** covariance structure
- Δ_k is the **separate column** covariance structure

2. **Estimate** covariances $\Sigma, \Delta_1, \dots, \Delta_K$ to obtain $\hat{\Sigma}, \hat{\Delta}_1, \dots, \hat{\Delta}_K$

$$\hat{\Sigma}, \hat{\Delta}_1, \dots, \hat{\Delta}_K = \operatorname{argmax} \ell(\Sigma, \Delta_1, \dots, \Delta_K)$$

3. **Maximize** covariances

iPC Scores: $\mathbf{U} \leftarrow$ eigenvectors of $\hat{\Sigma}$ = joint patterns

iPC Loadings: $\mathbf{V}_k \leftarrow$ eigenvectors of $\hat{\Delta}_k$ = individual patterns

4. **Visualize** dominant joint patterns by plotting iPC scores \mathbf{U}

1. iPCA Model: For each $k = 1, \dots, K$, assume

$$\mathbf{X}_k \sim N_{n,p_k}(\mathbf{1}_n \mu_k^T, \Sigma \otimes \Delta_k)$$

- Σ is the **common row** covariance structure
- Δ_k is the **separate column** covariance structure

2. Estimate covariances $\Sigma, \Delta_1, \dots, \Delta_K$ to obtain $\hat{\Sigma}, \hat{\Delta}_1, \dots, \hat{\Delta}_K$

$$\hat{\Sigma}, \hat{\Delta}_1, \dots, \hat{\Delta}_K = \operatorname{argmax} \ell(\Sigma, \Delta_1, \dots, \Delta_K) - \sum_{k=1}^K \lambda_k \|\Sigma^{-1} \otimes \Delta_k^{-1}\|_F^2$$

3. Maximize covariances

iPC Scores: $\mathbf{U} \leftarrow$ eigenvectors of $\hat{\Sigma}$ = joint patterns

iPC Loadings: $\mathbf{V}_k \leftarrow$ eigenvectors of $\hat{\Delta}_k$ = individual patterns

4. Visualize dominant joint patterns by plotting iPC scores \mathbf{U}

CASE STUDY: ALZHEIMER'S DISEASE (AD)

ROSMAP Study (*Bennett et al. (2012)*)

- Longitudinal clinical-pathological cohort study of aging and AD
- Genomics on post-mortem brains

ROSMAP Genomics Data ($n = 507$)

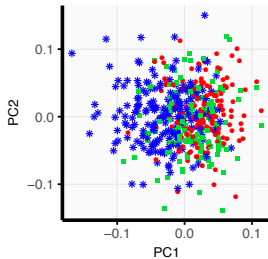
- miRNA Expression: $p_1 = 309$
- Gene Expression via RNASeq (log-transformed): $p_2 = 900$
- DNA Methylation (m-values): $p_3 = 1250$

Clinical Outcomes of Interest:

- Clinician's Diagnosis
 - AD, Mild Cognitive Impairment, No Cognitive Impairment
- Global Cognition Score

ROSMAP: CLINICIAN'S DIAGNOSIS

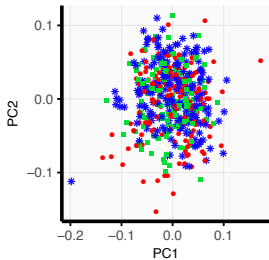
iPCA (x Frobenius)



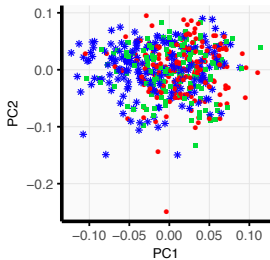
Diagnosis

- NCI
- MCI
- * AD

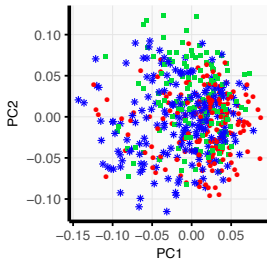
PCA on miRNA



PCA on RNASeq

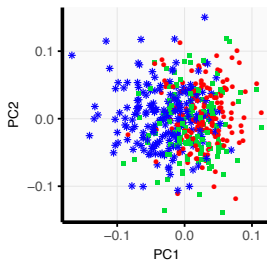


PCA on Methylation



ROSMAP: CLINICIAN'S DIAGNOSIS

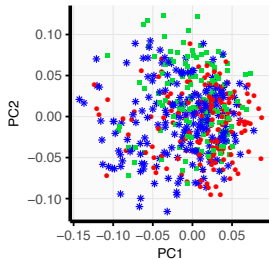
iPCA (x Frobenius)



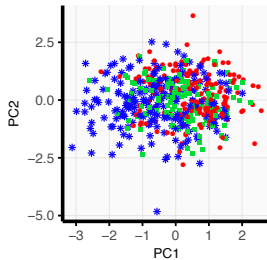
Diagnosis

- NCI
- MCI
- * AD

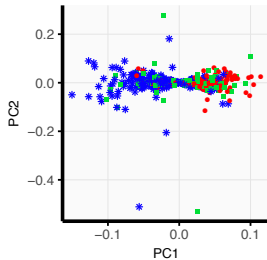
Concatenated



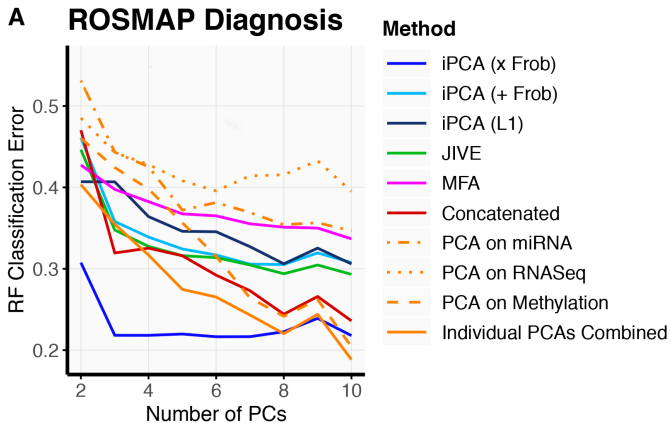
MFA



JIVE

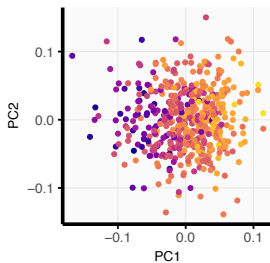


ROSMAP: CLINICIAN'S DIAGNOSIS

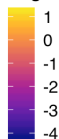


ROSMAP: GLOBAL COGNITION

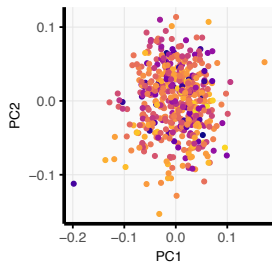
iPCA (x Frobenius)



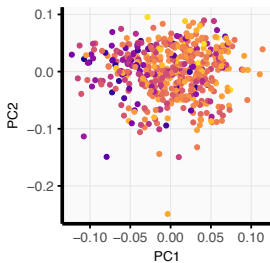
Cognition



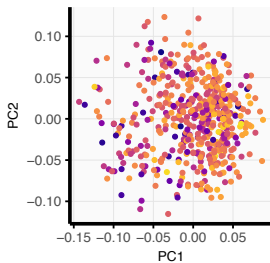
PCA on miRNA



PCA on RNASeq

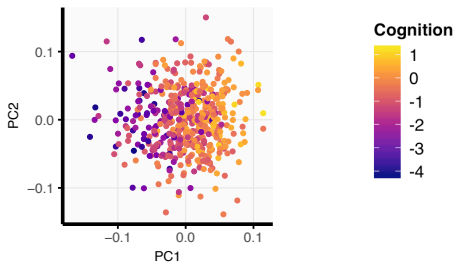


PCA on Methylation

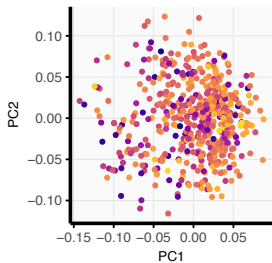


ROSMAP: GLOBAL COGNITION

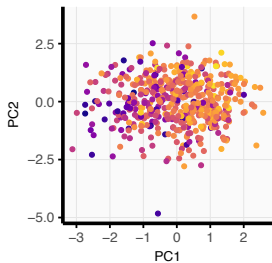
iPCA (x Frobenius)



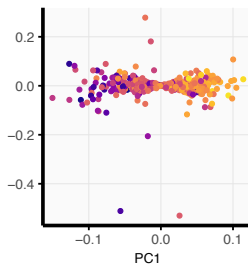
Concatenated



MFA



JIVE



B ROSMAP Cognition

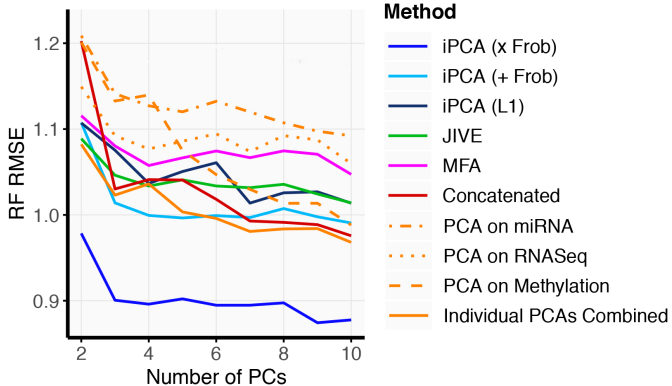


Table 1: Top genetic features obtained by applying Sparse PCA to each $\hat{\Delta}_k$ in ROSMAP analysis (using the multiplicative Frobenius iPCA estimator)

	miRNA	RNASeq	Methylation
1	miR 216a	VGF	TMCO6
2	miR 127 3p	SVOP	PHF3
3	miR 124	PCDHGC5	BRUNOL4
4	miR 30c	ADCYAP1	OSCP1
5	miR 143	LINC01007	GRIN2B

CONCLUSION

- iPCA is a new practical tool for discovering and visualizing interesting joint patterns which occur in multiple datasets
- In order to fit iPCA, we propose using the multiplicative Frobenius estimator.
- Though we impose a model, the assumptions are analogous to those in classical PCA
- Tang, T. M., & Allen, G. I. (2018). Integrated Principal Components Analysis. *arXiv preprint arXiv:1810.00832*.

ACKNOWLEDGMENTS

Collaborator

- **Genevera Allen**

Thank you also to...

- ROSMAP Team
 - Joshua Shulman
 - David Bennett
 - Zhandong Liu
 - Ying-Wooi Wan
- Bin Yu

Funding

- National Science Foundation
Graduate Research Fellowship
- Office of Naval Research



Thank you!