

# Bidimensional Linked Matrix Decomposition for Pan-Omics Pan-Cancer Analysis

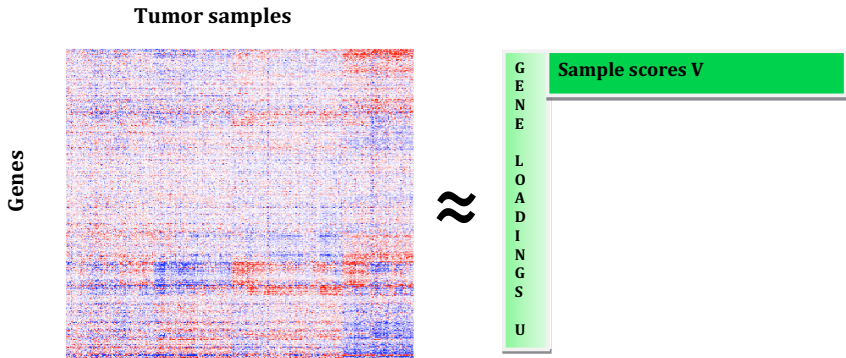
Eric F. Lock

University of Minnesota, Division of Biostatistics

JSM Virtual, 08/06/2020

# Matrix factorization

- Gene expression matrix  $X : m \times n$ 
  - $m$  genes for  $n$  breast cancer tumor samples

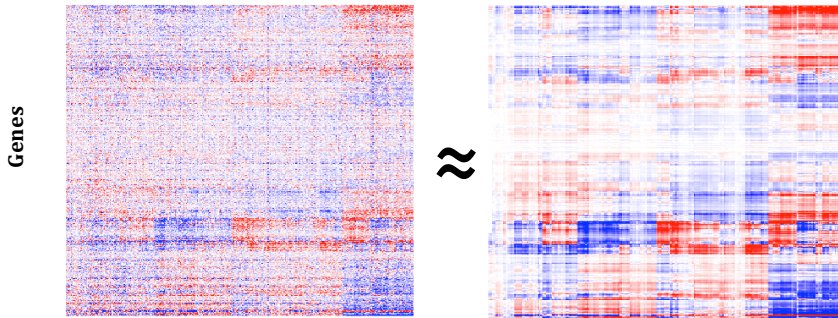


- Low rank factorization:  $X \approx UV$ ,  $U : m \times r$ ,  $V : r \times n$ .

# Matrix factorization ( $r=3$ )

- Gene expression matrix  $X : m \times n$ 
  - $m$  genes for  $n$  breast cancer tumor samples

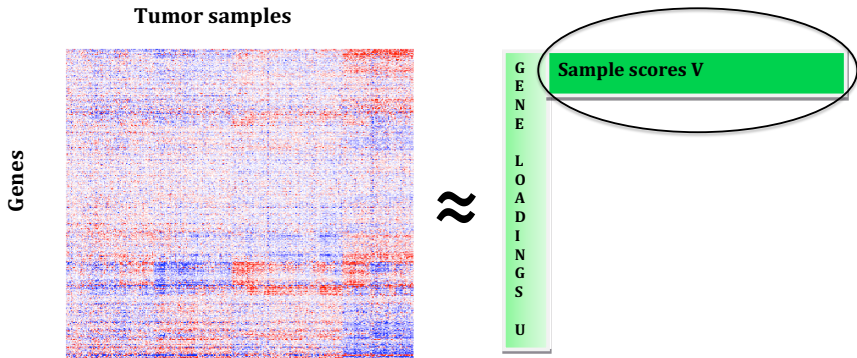
**Tumor samples**



- Low rank factorization:  $X \approx UV$ ,  $U : m \times 3$ ,  $V : 3 \times n$ .

# Matrix factorization ( $r=3$ )

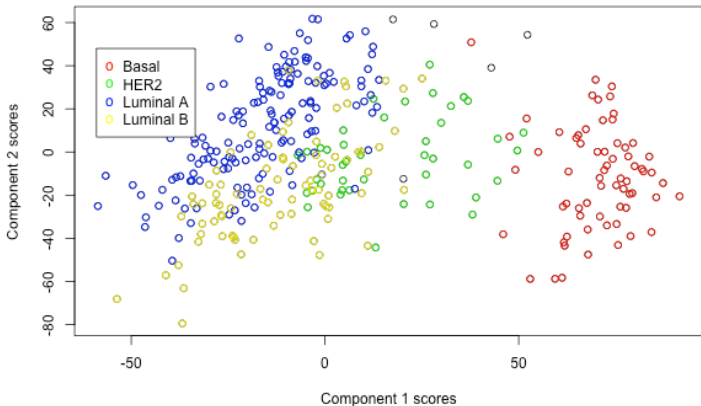
- Gene expression matrix  $X : m \times n$ 
  - $m$  genes for  $n$  breast cancer tumor samples



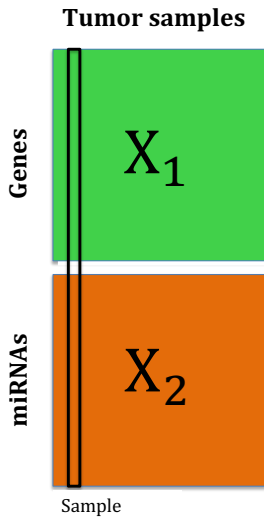
- Low rank factorization:  $X \approx UV$ ,  $U : m \times 3$ ,  $V : 3 \times n$ .

# Matrix factorization

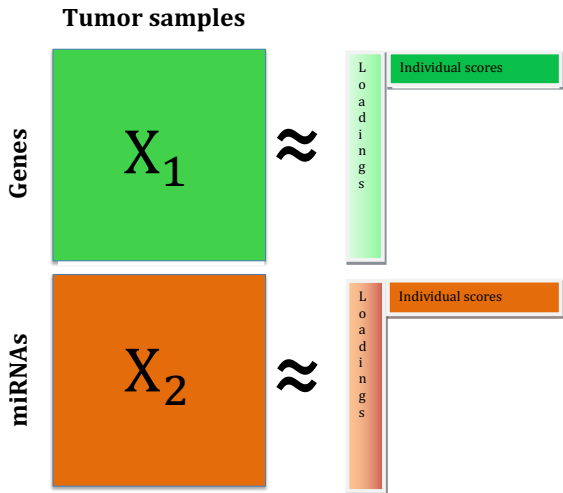
- First two principal component scores
  - Colored by breast tumor subtype



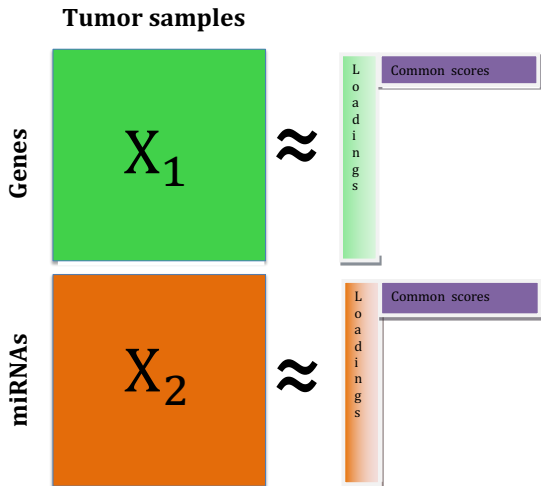
# Vertically linked data



# Vertically linked data: separate factorizations

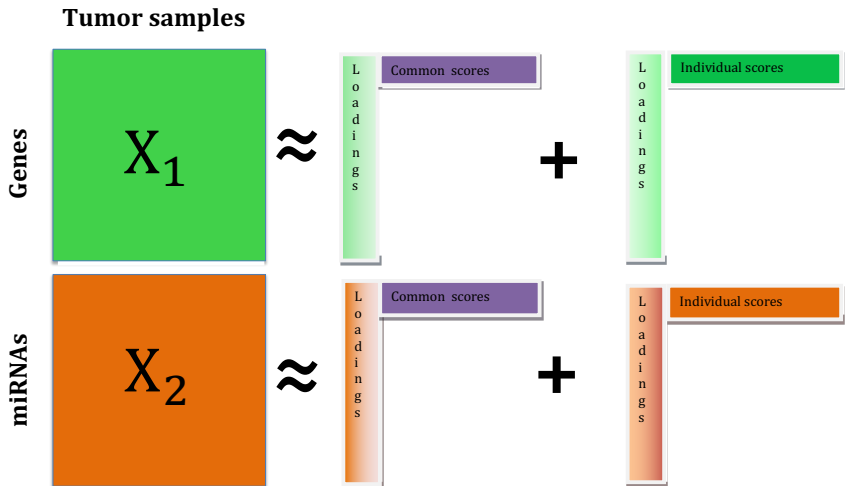


# Vertically linked data: joint factorization





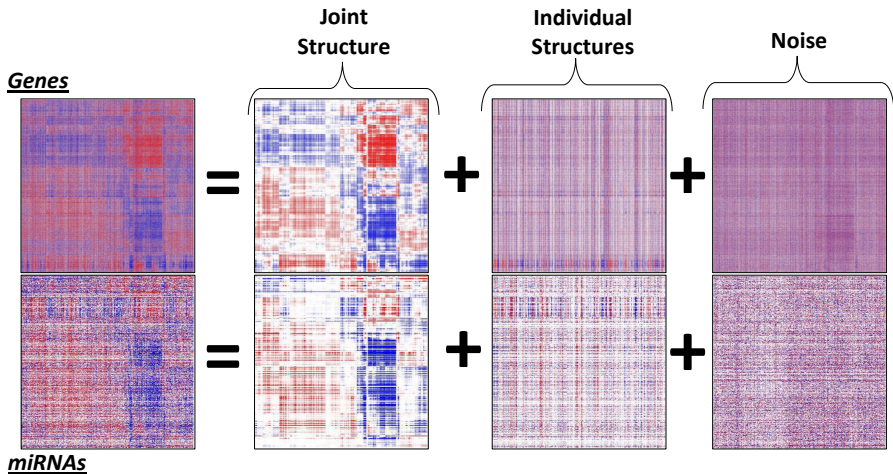
# Vertically linked data: JIVE factorization



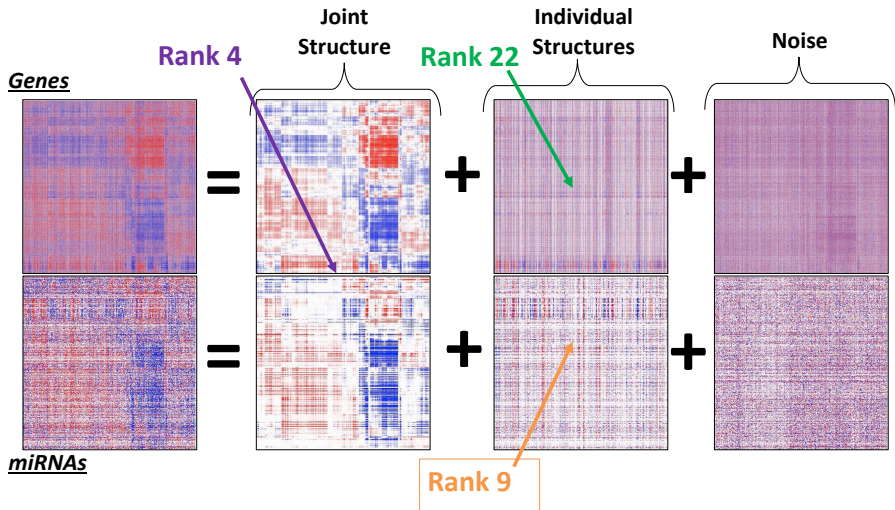
# Joint + individual factorization methods

- ▶ JIVE [Lock, Hoadley, Marron, and Nobel, 2013]
  - ▶ “Joint and Individual Variation Explained”
- ▶ R.JIVE [O’Connell and Lock, 2016]
- ▶ AJIVE [Feng, Jiang, Hannig and Marron, 2018]
- ▶ SLIDE [Gaynanova and Li, 2018]
- ▶ GIPCA [Zhu, Li, Lock, 2018]
- ▶ COBE, SIFA, MOFA, & more!

# JIVE Estimates

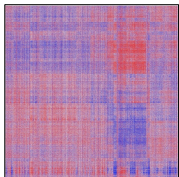


# JIVE Estimates

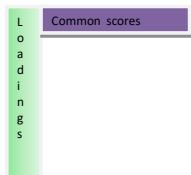


# JIVE Estimates (factorized)

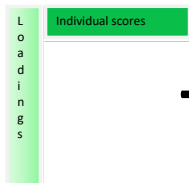
Genes



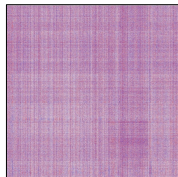
=



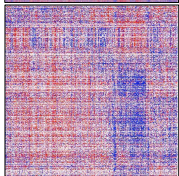
+



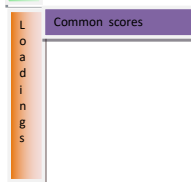
+



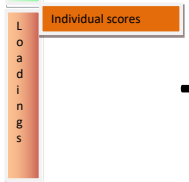
miRNAs



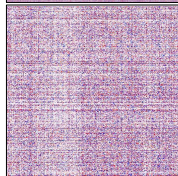
=



+

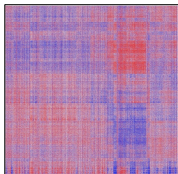


+

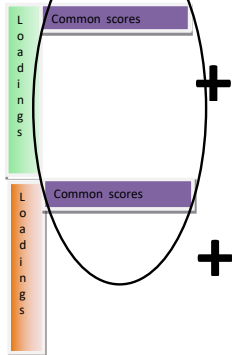


# JIVE Estimates (factorized)

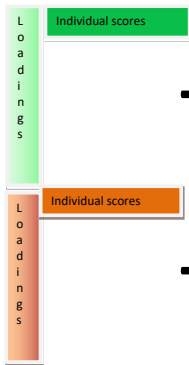
Genes



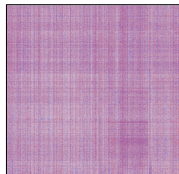
=



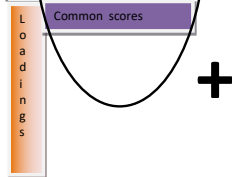
+



+



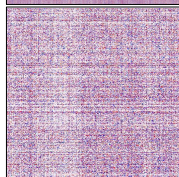
=



+

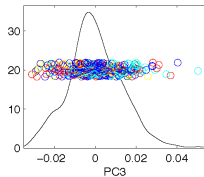
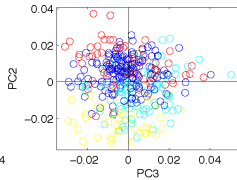
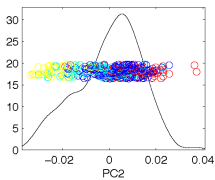
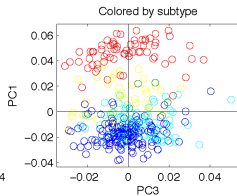
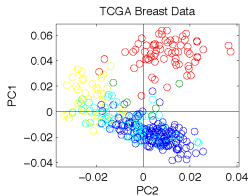
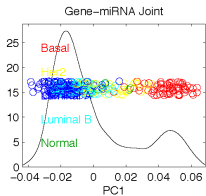


+

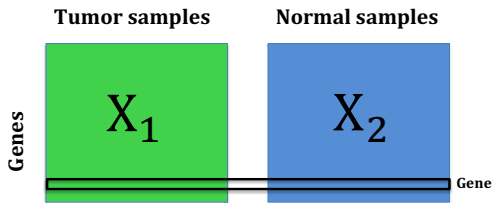


miRNAs

# Joint PCs

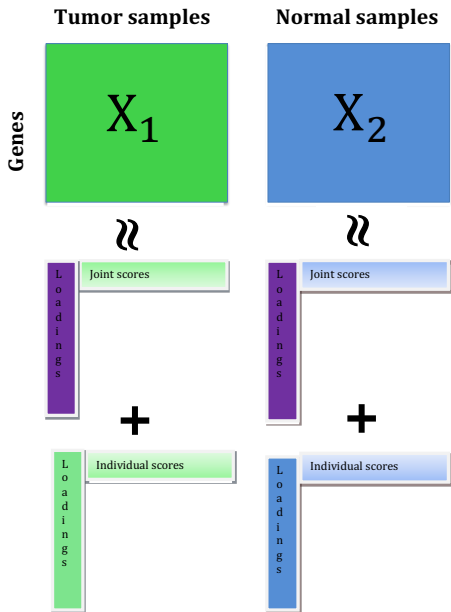


# Horizontally linked data

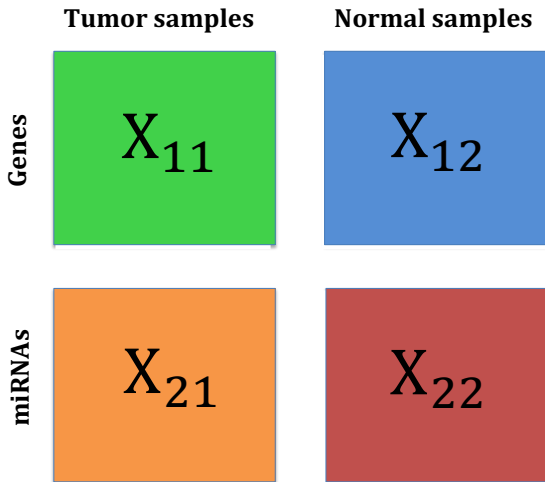




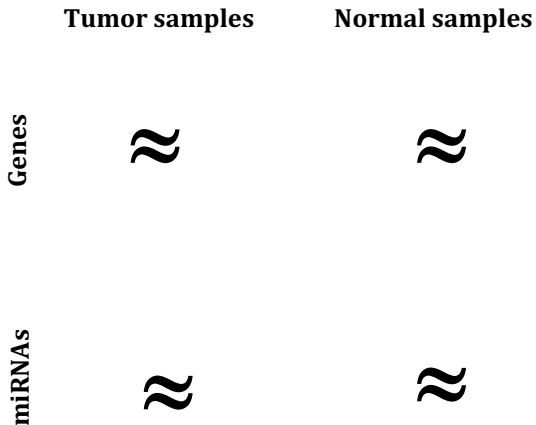
# Horizontally linked data: JIVE factorization



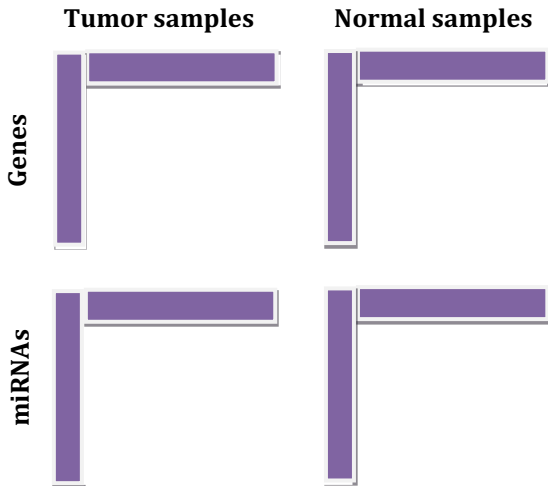
# Bidimensionally linked data: BIDIFAC



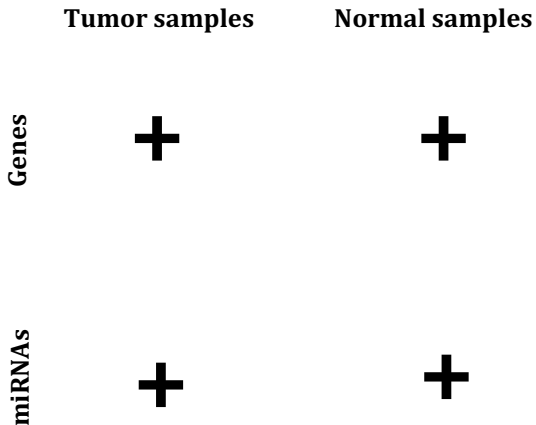
# Bidimensionally linked data: BIDIFAC



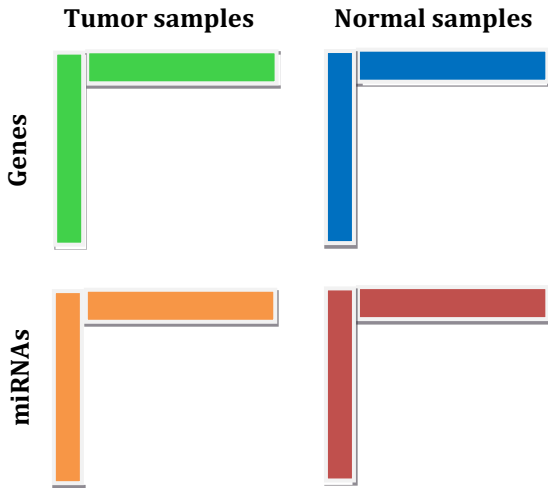
# Bidimensionally linked data: BIDIFAC



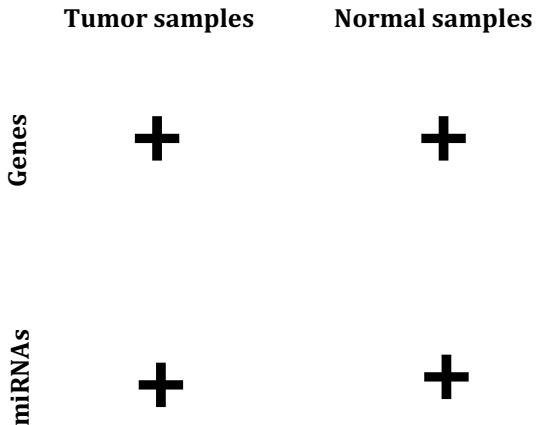
# Bidimensionally linked data: BIDIFAC



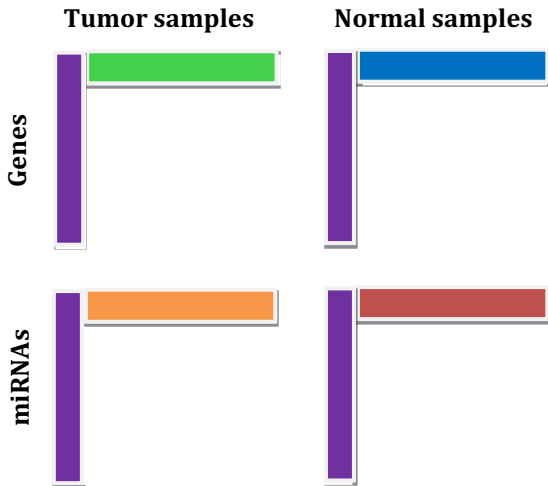
# Bidimensionally linked data: BIDIFAC



# Bidimensionally linked data: BIDIFAC

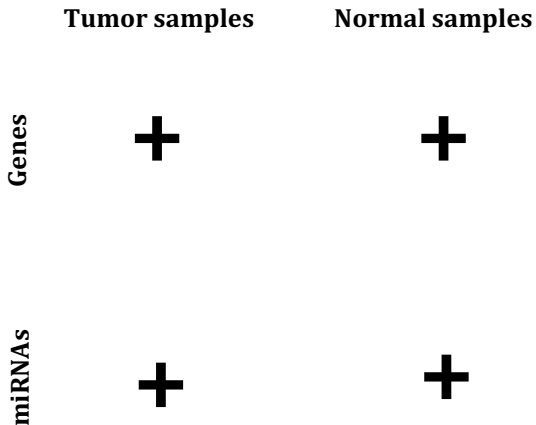


# Bidimensionally linked data: BIDIFAC

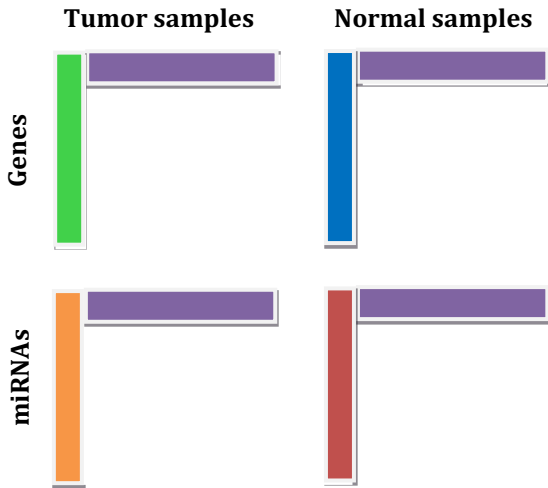




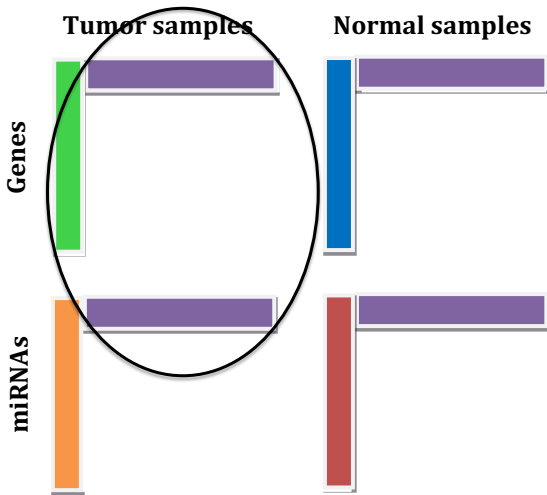
# Bidimensionally linked data: BIDIFAC



# Bidimensionally linked data: BIDIFAC



# Bidimensionally linked data: BIDIFAC



# Tumor-Specific Columns Shared PCs

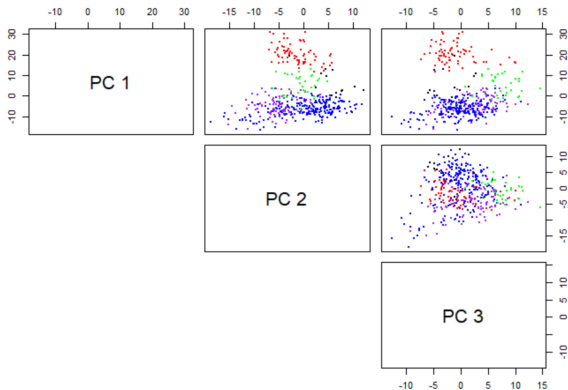
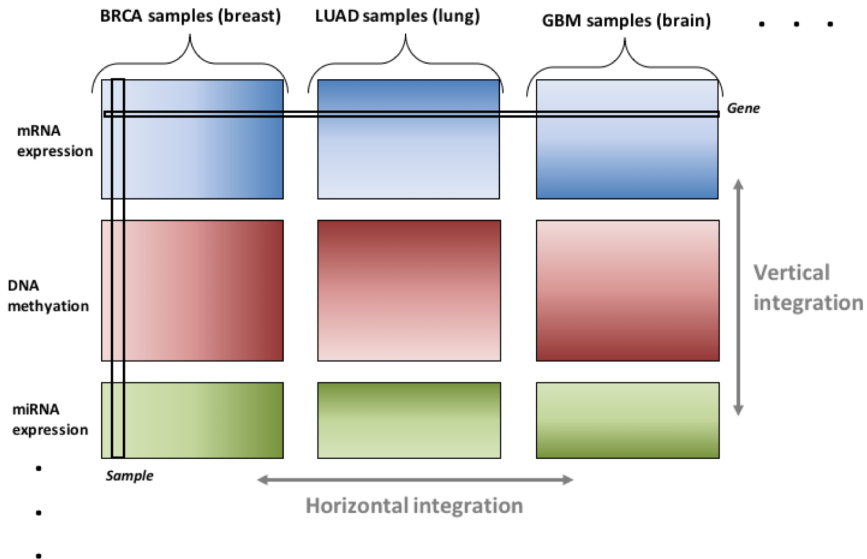


Figure: Principal components of the estimated column-shared structure, colored by subtype: Basal, HER2, Lum A, Lum B.

# Pan-omics pan-cancer integration!



# BIDIFAC: general framework

Consider matrices  $\{\mathbf{X}_{ij} : M_i \times N_j \mid i = 1, \dots, I, j = 1, \dots, J\}$ .

$$\mathbf{X}_{..} = \begin{bmatrix} \mathbf{X}_{11} & \dots & \mathbf{X}_{1I} \\ \vdots & \ddots & \vdots \\ \mathbf{X}_{J1} & \dots & \mathbf{X}_{IJ} \end{bmatrix}$$

$$\mathbf{X}_{.i} = [\mathbf{X}_{i1}, \dots, \mathbf{X}_{iq}]$$

$$\mathbf{X}_{.j} = \begin{bmatrix} \mathbf{X}_{1j} \\ \vdots \\ \mathbf{X}_{pj} \end{bmatrix}$$

- Decompose  $\mathbf{X}_{..}$  into low-rank structural *modules*:

$$\mathbf{X}_{..} = \sum_{k=1}^{\kappa} \mathbf{S}_{..}^{(k)} + \mathbf{E}_{..}, \quad (1)$$

where

$$\mathbf{S}_{00}^{(k)} = \begin{bmatrix} \mathbf{S}_{11}^{(k)} & \mathbf{S}_{12}^{(k)} & \cdots & \mathbf{S}_{1q}^{(k)} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{S}_{p1}^{(k)} & \mathbf{S}_{p2}^{(k)} & \cdots & \mathbf{S}_{pq}^{(k)} \end{bmatrix}$$

presence of each  $\mathbf{S}_{ij}^{(k)}$  is determined by  $\mathbf{R} : I \times \kappa$  and  $\mathbf{C} : J \times \kappa$ :

$$\mathbf{S}_{ij}^{(k)} = \begin{cases} \mathbf{0}_{M_i \times N_j} & \text{if } \mathbf{R}[i, k] = 0 \text{ or } \mathbf{C}[j, k] = 0 \\ \mathbf{U}_i^{(k)} \mathbf{V}_j^{(k)} & \text{if } \mathbf{R}[i, k] = 1 \text{ and } \mathbf{C}[j, k] = 1 \end{cases}.$$

- ▶ Minimize the following objective over  $\mathbf{R}$ ,  $\mathbf{C}$ , and  $\{\mathbf{S}_{..}^{(k)}\}_{k=1}^{\kappa}$ :

$$\|\mathbf{X}_{..} - \sum_{k=1}^{\kappa} \mathbf{S}_{..}^{(k)}\|_F^2 + \sum_{k=1}^{\kappa} \lambda_k \|\mathbf{S}_{..}^{(k)}\|_*$$

where  $\|\cdot\|_*$  gives the nuclear norm

$$\text{SVD}(\mathbf{A}) = \mathbf{U}\mathbf{D}\mathbf{V}^T \rightarrow \|\mathbf{A}\|_* = \sum \mathbf{D}[i, i].$$

- ▶ Choice of  $\lambda_k$ 's
  - ▶ Determined by random matrix theory for singular values  $\mathbf{D}$
  - ▶ Gaurantees each module  $\mathbf{S}_{..}^{(k)}$  is low-rank.
  - ▶ Any submatrix  $(\mathbf{R}[, k], \mathbf{C}[, k])$  can have a non-zero module



- ▶ Let  $\mathbb{S}_{\hat{\mathbf{X}}}$  be the set of possible decompositions for  $\hat{\mathbf{X}}_{..}$ :

$$\mathbb{S}_{\hat{\mathbf{X}}} = \left\{ \{ \mathbf{s}_{..}^{(k)} \}_{k=1}^K \mid \hat{\mathbf{X}}_{..} = \sum_{k=1}^K \mathbf{s}_{..}^{(k)} \right\}.$$

## Theorem

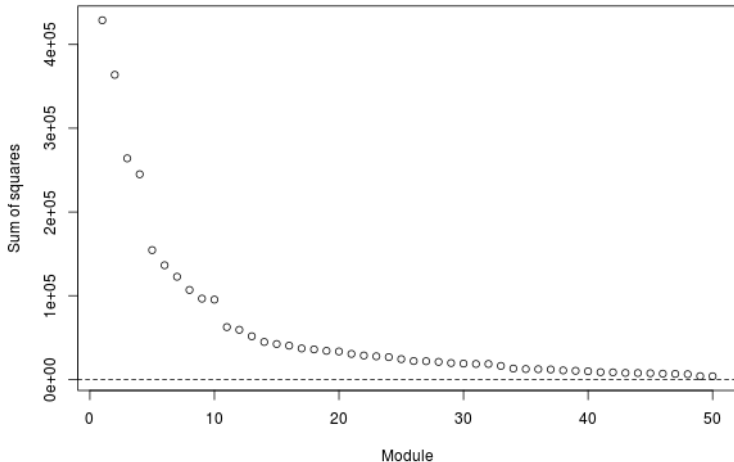
Consider  $\{ \hat{\mathbf{S}}_{..}^{(k)} \}_{k=1}^K \in \mathbb{S}_{\hat{\mathbf{X}}}$  and let  $\mathbf{U}_{.}^{(k)} \hat{\mathbf{D}}^{(k)} \mathbf{V}_{.}^{(k)T}$  give the SVD of  $\hat{\mathbf{S}}_{..}^{(k)}$ . The following three properties uniquely identify  $\{ \hat{\mathbf{S}}_{..}^{(k)} \}_{k=1}^K$ .

- ▶  $\{ \hat{\mathbf{S}}_{..}^{(k)} \}_{k=1}^K$  minimizes  $\sum_{k=1}^K \lambda_k \| \mathbf{s}_{..}^{(k)} \|_*$  over  $\mathbb{S}_{\hat{\mathbf{X}}}$ ,
- ▶  $\{ \hat{\mathbf{U}}_i^{(k)}[\cdot, r] : \mathbf{R}[i, k] = 1 \text{ and } \hat{\mathbf{D}}^{(k)}[r, r] > 0 \}$  are linearly independent for  $i = 1, \dots, p$ ,
- ▶  $\{ \hat{\mathbf{V}}_j^{(k)}[\cdot, r] : \mathbf{C}[j, k] = 1 \text{ and } \hat{\mathbf{D}}^{(k)}[r, r] > 0 \}$  are linearly independent for  $j = 1, \dots, q$ .

- ▶ TCGA data for 6793 samples representing 29 cancer types:
  - ▶ ACC, BLCA, BRCA, CESC, CHOL, CORE, DLBC, ESCA, HNSC, KICH, KIRC, KIRP, LGG, LIHC, LUAD, LUSC, MESO, OV, PAAD, PCPG, PRAD, SARC, SKCM, STAD, TGCT, THCA, THYM, UCEC, and UCS.
- ▶ Data for 4 different 'omics platforms
  - ▶ Gene expression (mRNA), miRNA, DNA methylation, and protein abundance
- ▶  $(2^4 - 1) \cdot (2^{29} - 1) = 8053063665$  possible modules!

# BIDIFAC+: results

- Variance explained for each module  $\mathbf{S}_{..}^{(k)}$  ( $K = 50$ ):



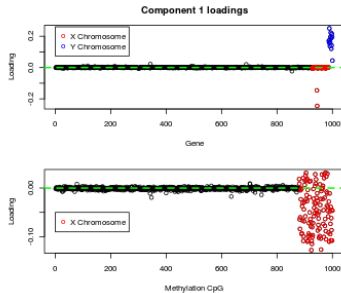
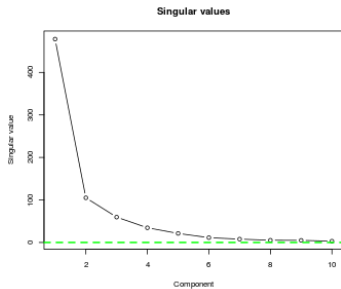
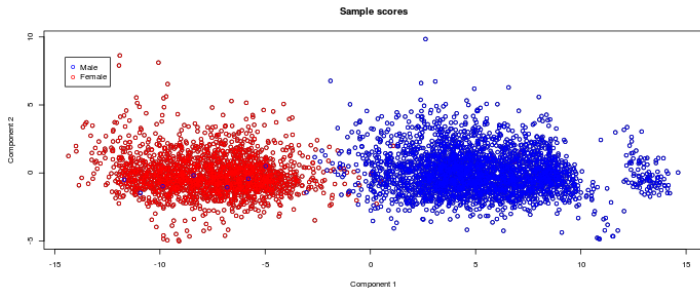
- Top structural modules, ranked by variance explained:

Module	Cancer types	Omics sources
1	All cancers	mRNA miRNA Meth Protein
2	All cancers	miRNA
3	BLCA BRCA CESC CHOL CORE DLBC ESCA HNSC LIHC LUAD LUSC OV PAAD PRAD SKCM STAD TGCT UCEC UCS	Meth
4	ACC BLCA CHOL CORE DLBC ESCA HNSC KICH KIRC KIRP LGG LIHC LUAD LUSC MESO PAAD PCPG SARC SKCM STAD THCA THYM	mRNA Meth
5	All cancers	mRNA
6	BRCA	mRNA miRNA Meth Protein
7	LGG	mRNA miRNA Protein
8	All cancers *but* LGG	Protein
9	THCA	mRNA miRNA Protein
10	All cancers *but* LGG and TGCT	miRNA
11	CHOL KIRC KIRP LIHC	mRNA miRNA Meth Protein
12	LGG	Meth
13	BLCA CESC CORE ESCA HNSC LUSC SARC STAD	mRNA miRNA Meth Protein
14	KICH KIRC KIRP	mRNA miRNA Protein
15	BLCA BRCA CESC CHOL ESCA HNSC LUAD LUSC PAAD PRAD SKCM STAD TGCT UCEC UCS	mRNA miRNA

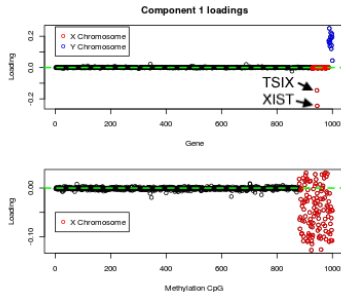
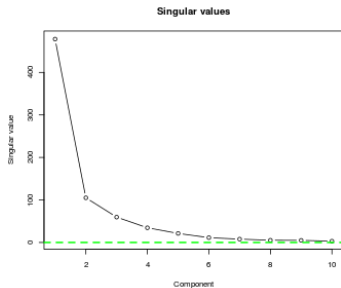
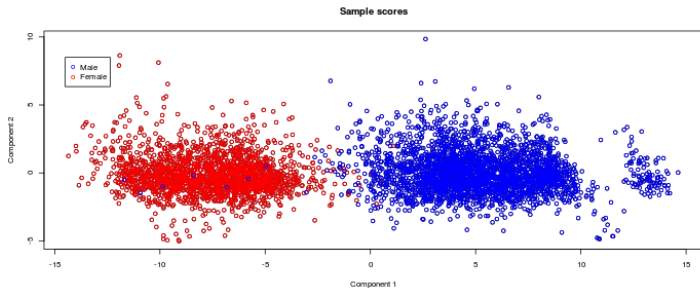
- Top structural modules, ranked by variance explained:

Module	Cancer types	Omics sources
1	All cancers	mRNA miRNA Meth Protein
2	All cancers	miRNA
3	BLCA BRCA CESC CHOL CORE DLBC ESCA HNSC LIHC LUAD LUSC OV PAAD PRAD SKCM STAD TGCT UCEC UCS	Meth
4	ACC BLCA CHOL CORE DLBC ESCA HNSC KICH KIRC KIRP LGG LIHC LUAD LUSC MESO PAAD PCPG SARC SKCM STAD THCA THYM	mRNA Meth
5	All cancers	mRNA
6	BRCA	mRNA miRNA Meth Protein
7	LGG	mRNA miRNA Protein
8	All cancers *but* LGG	Protein
9	THCA	mRNA miRNA Protein
10	All cancers *but* LGG and TGCT	miRNA
11	CHOL KIRC KIRP LIHC	mRNA miRNA Meth Protein
12	LGG	Meth
13	BLCA CESC CORE ESCA HNSC LUSC SARC STAD	mRNA miRNA Meth Protein
14	KICH KIRC KIRP	mRNA miRNA Protein
15	BLCA BRCA CESC CHOL ESCA HNSC LUAD LUSC PAAD PRAD SKCM STAD TGCT UCEC UCS	mRNA miRNA

# Module 4



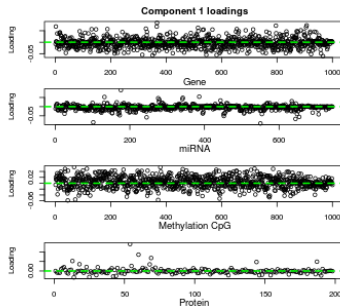
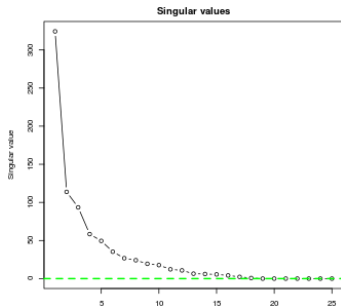
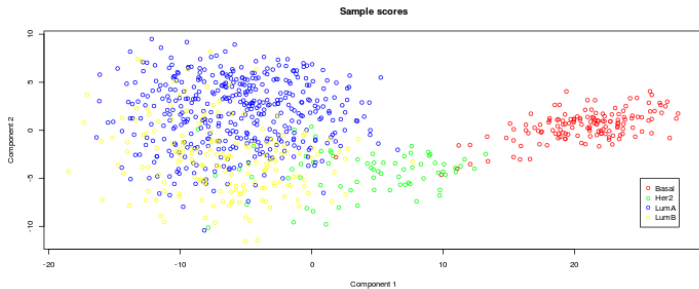
# Module 4



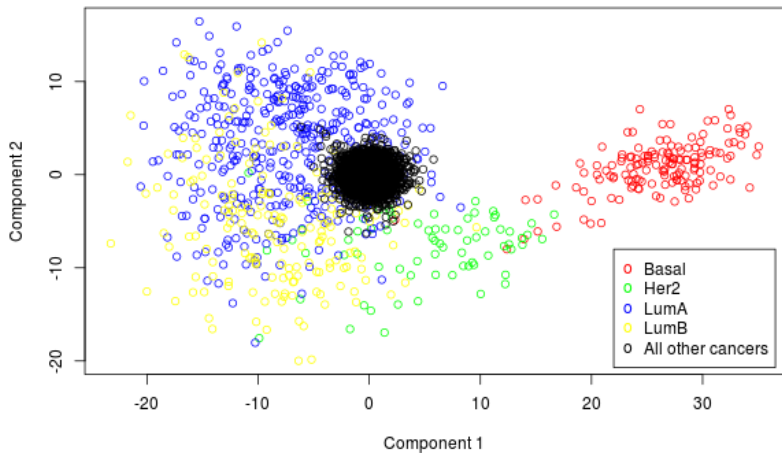
- Top structural modules, ranked by variance explained:

Module	Cancer types	Omics sources
1	All cancers	mRNA miRNA Meth Protein
2	All cancers	miRNA
3	BLCA BRCA CESC CHOL CORE DLBC ESCA HNSC LIHC LUAD LUSC OV PAAD PRAD SKCM STAD TGCT UCEC UCS	Meth
4	ACC BLCA CHOL CORE DLBC ESCA HNSC KICH KIRC KIRP LGG LIHC LUAD LUSC MESO PAAD PCPG SARC SKCM STAD THCA THYM	mRNA Meth
5	All cancers	mRNA
6	<b>BRCA</b>	<b>mRNA miRNA Meth Protein</b>
7	LGG	mRNA miRNA Protein
8	All cancers *but* LGG	Protein
9	THCA	mRNA miRNA Protein
10	All cancers *but* LGG and TGCT	miRNA
11	CHOL KIRC KIRP LIHC	mRNA miRNA Meth Protein
12	LGG	Meth
13	BLCA CESC CORE ESCA HNSC LUSC SARC STAD	mRNA miRNA Meth Protein
14	KICH KIRC KIRP	mRNA miRNA Protein
15	BLCA BRCA CESC CHOL ESCA HNSC LUAD LUSC PAAD PRAD SKCM STAD TGCT UCEC UCS	mRNA miRNA





Sample scores



# Thank you!

- ▶ Support: NCI grant R21CA231214-01
- ▶ References:
  - ▶ **BIDIFAC**: J Park & EF Lock. Integrative Factorization of Bidimensionally Linked Matrices. *Biometrics*, 76 (1): 61-74 2020.
  - ▶ **BIDIFAC+**: EF Lock, J Park & KA Hoadley. Bidimensional linked matrix factorization for pan-omics pan-cancer analysis. *Preprint*, arXiv:2002.0260, 2020.
- ▶ Code:
  - ▶ BIDIFAC: <https://github.com/lockEF/bidifac>