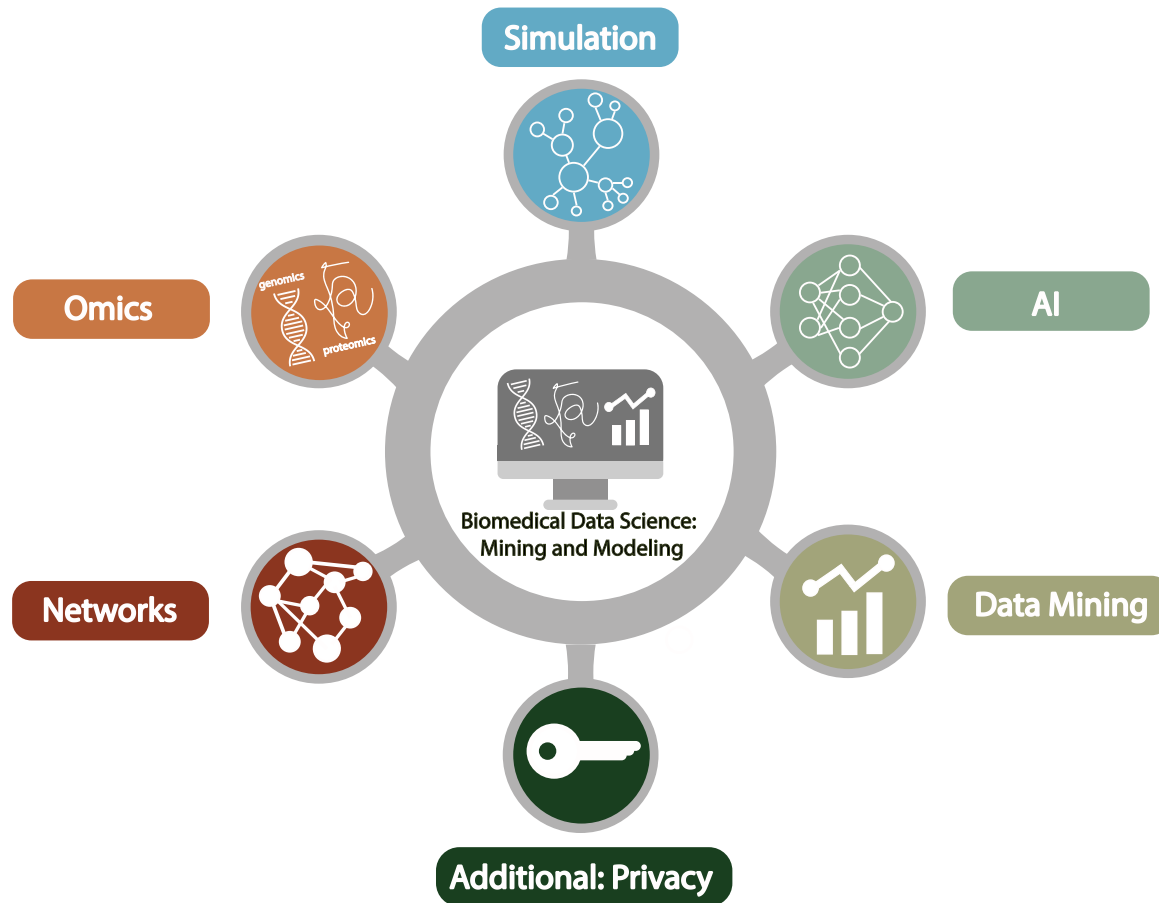


Biomedical Data Science (GersteinLab.org/courses/452)

Unsupervised Datamining – SVD (23m9c)



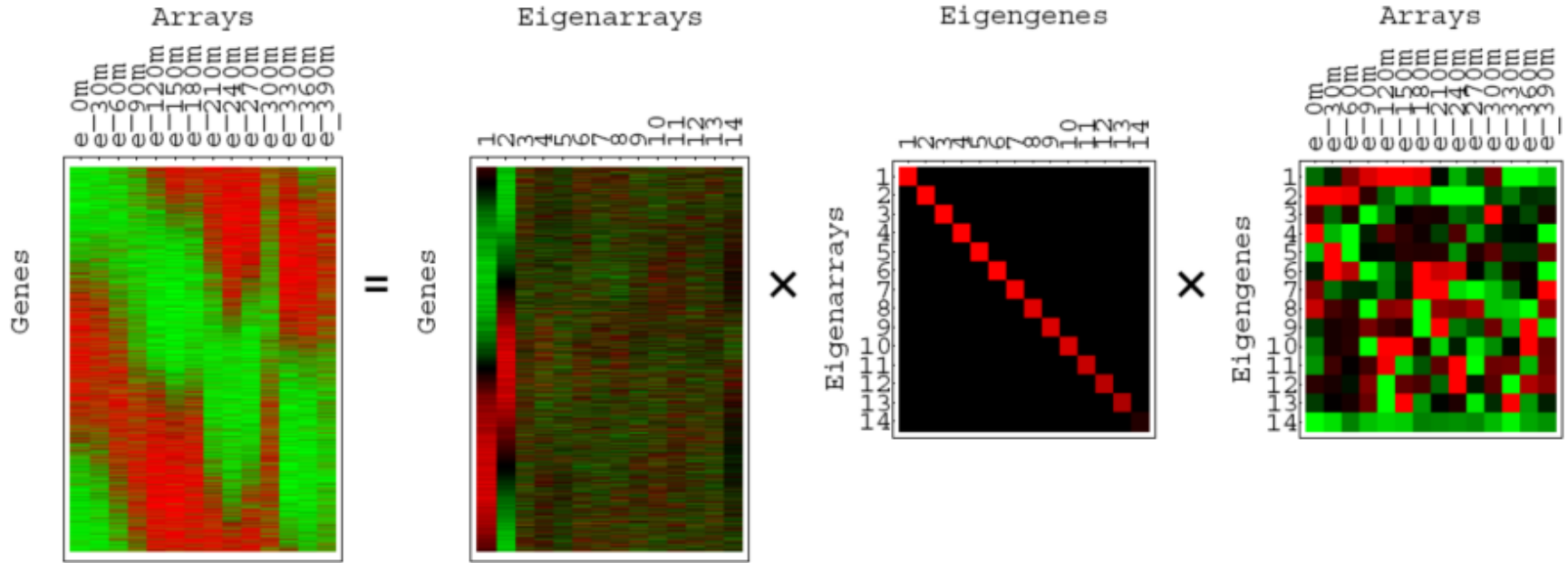
Last edit in spring '23. Condensing by
~3 slide deletions 2022's 22m9c,
which is similar to
2021's M9c [which has a video].

Unsupervised Mining

SVD

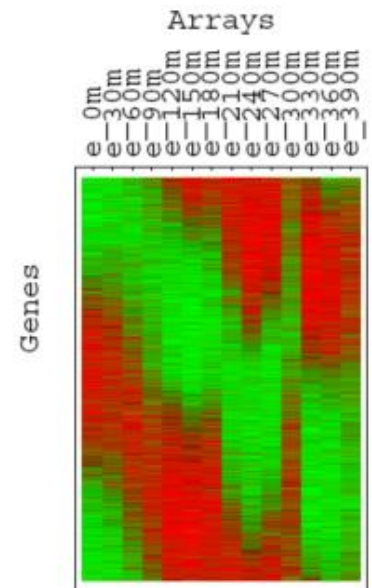
Puts together slides prepared by
Brandon Xia with images from
Alter et al. papers

SVD for microarray data (Alter et al, PNAS 2000)



$$A = USV^T$$

- A is any rectangular matrix ($m \geq n$)
- Row space: vector subspace generated by the row vectors of A
- Column space: vector subspace generated by the column vectors of A
 - The dimension of the row & column space is the rank of the matrix A : $r (\leq n)$
- A is a linear transformation that maps vector x in row space into vector Ax in column space

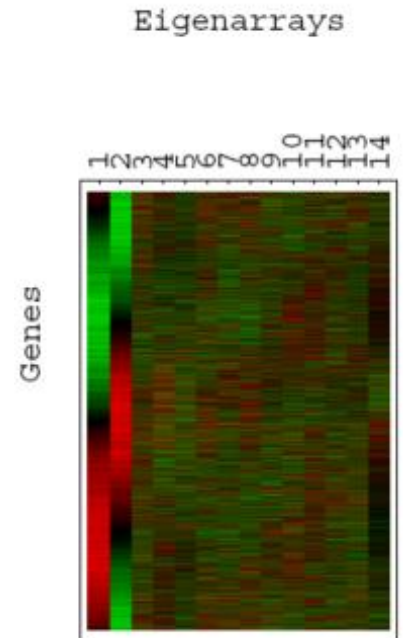


$$A = USV^T$$

- U is an “orthogonal” matrix ($m \geq n$)
- Column vectors of U form an orthonormal basis for the **column space** of A: $U^T U = I$

$$U = \begin{pmatrix} | & | & & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_n \\ | & | & & | \end{pmatrix}$$

- $\mathbf{u}_1, \dots, \mathbf{u}_n$ in U are eigenvectors of AA^T
 - $AA^T = USV^T VSU^T = US^2 U^T$
 - “Left singular vectors”

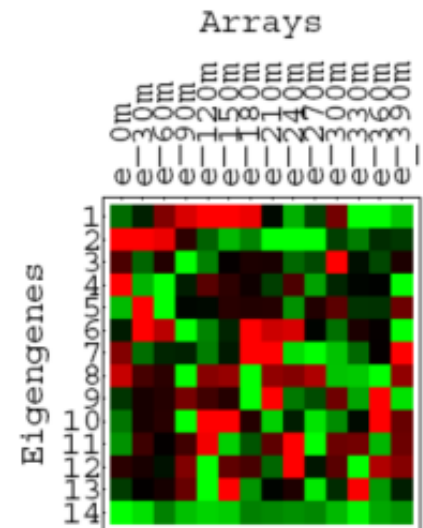


$$A = USV^T$$

- V is an orthogonal matrix (n by n)
- Column vectors of V form an orthonormal basis for the **row space** of A : $V^T V = V V^T = I$

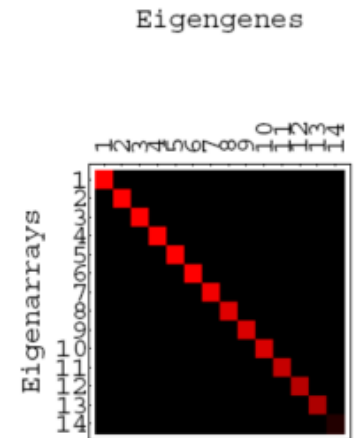
$$V = \begin{pmatrix} | & | & & | \\ \mathbf{v}_1 & \mathbf{v}_2 & \dots & \mathbf{v}_n \\ | & | & & | \end{pmatrix}$$

- $\mathbf{v}_1, \dots, \mathbf{v}_n$ in V are eigenvectors of $A^T A$
 - $A^T A = V S U^T U S V^T = V S^2 V^T$
 - “Right singular vectors”



$$A = USV^T$$

- S is a diagonal matrix (n by n) of non-negative singular values
- Typically sorted from largest to smallest
- Singular values are the non-negative square root of corresponding eigenvalues of $A^T A$ and AA^T



$$AV = US$$

- Means each $A\mathbf{v}_i = s_i\mathbf{u}_i$
- Remember A is a linear map from row space to column space
- Here, A maps an orthonormal basis $\{\mathbf{v}_i\}$ in row space into an orthonormal basis $\{\mathbf{u}_i\}$ in column space
- Each component of \mathbf{u}_i is the projection of a row of the data matrix A onto the vector \mathbf{v}_i

SVD as sum of rank-1 matrices

- $A = USV^T$
- $A = s_1 \mathbf{u}_1 \mathbf{v}_1^T + s_2 \mathbf{u}_2 \mathbf{v}_2^T + \dots + s_n \mathbf{u}_n \mathbf{v}_n^T$
- $s_1 \geq s_2 \geq \dots \geq s_n \geq 0$

an outer product
(uv^T) giving a
matrix rather than
the scalar of the
inner product

- What is the rank- r matrix \hat{A} that best approximates A ?

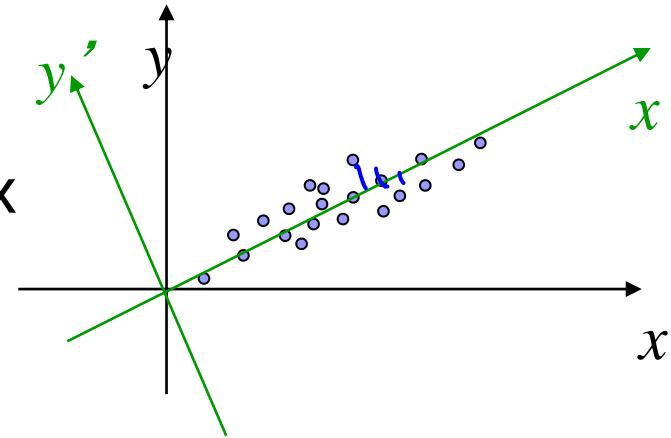
– Minimize
$$\sum_{i=1}^m \sum_{j=1}^n (\hat{A}_{ij} - A_{ij})^2$$

LSQ approx. If $r=1$,
this amounts to a
line fit.

- $\hat{A} = s_1 \mathbf{u}_1 \mathbf{v}_1^T + s_2 \mathbf{u}_2 \mathbf{v}_2^T + \dots + s_r \mathbf{u}_r \mathbf{v}_r^T$
- Very useful for matrix approximation

Geometry of SVD in row space

- A as a collection of m row vectors (points) in the row space of A
- $s_1 \mathbf{u}_1 \mathbf{v}_1^T + s_2 \mathbf{u}_2 \mathbf{v}_2^T$ is the best rank-2 matrix approximation for A
- Geometrically: \mathbf{v}_1 and \mathbf{v}_2 are the directions of the best approximating rank-2 subspace that goes through origin
- $s_1 \mathbf{u}_1$ and $s_2 \mathbf{u}_2$ gives coordinates for row vectors in rank-2 subspace
- \mathbf{v}_1 and \mathbf{v}_2 gives coordinates for row space basis vectors in rank-2 subspace



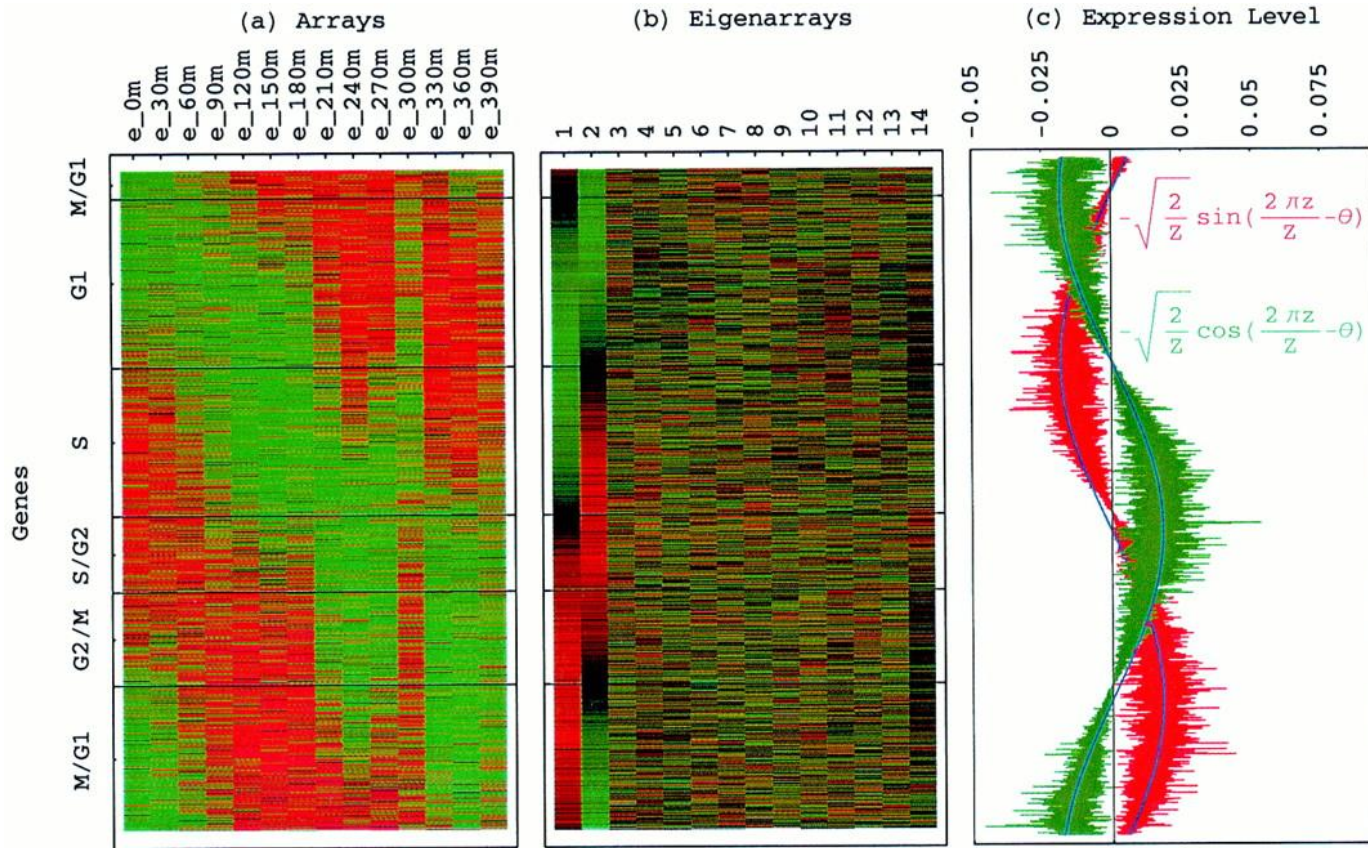
$$A \mathbf{v}_i = s_i \mathbf{u}_i$$

$$I \mathbf{v}_i = \mathbf{v}_i$$

Unsupervised Mining

Intuition on interpretation of SVD
in terms of genes and conditions

Genes sorted by correlation with top 2 eigengenes



Alter, Orly et al. (2000) Proc. Natl. Acad. Sci. USA 97, 10101-10106

Fig. 3. Genes sorted by relative correlation with $|\gamma_1\rangle_N$ and $|\gamma_2\rangle_N$ of normalized elutriation. (a) Normalized elutriation expression of the sorted 5,981 genes in the 14 arrays, showing traveling wave of expression. (b) Eigenarrays expression; the expression of $|\alpha_1\rangle_N$ and $|\alpha_2\rangle_N$, the eigenarrays corresponding to $|\gamma_1\rangle_N$ and $|\gamma_2\rangle_N$, displays the sorting. (c) Expression levels of $|\alpha_1\rangle_N$ (red) and $|\alpha_2\rangle_N$ (green) fit normalized sine and cosine functions of period $Z = N - 1 = 5,980$ and phase $\theta \approx 2\pi/13$ (blue), respectively.

Normalized elutriation expression in the subspace associated with the cell cycle

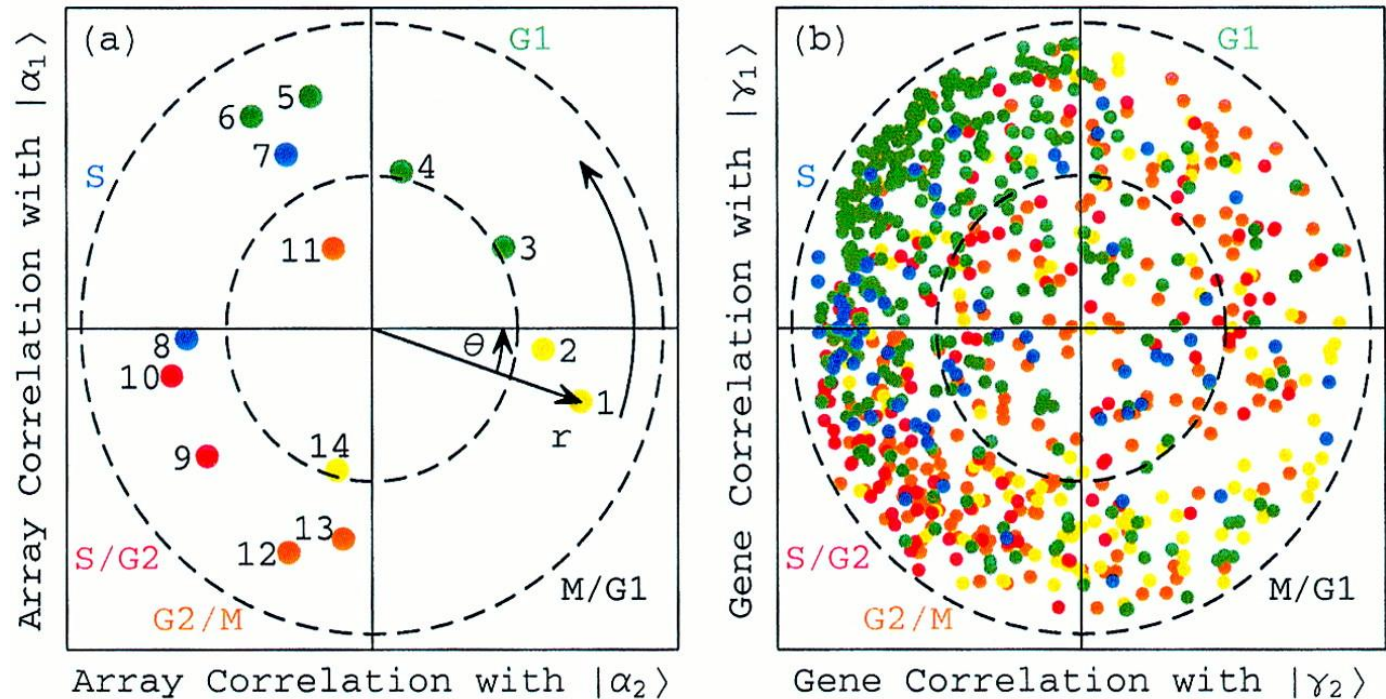


Fig. 2. Normalized elutriation expression in the subspace associated with the cell cycle. (a) Array correlation with $|\alpha_1\rangle_N$ along the y-axis vs. that with $|\alpha_2\rangle_N$ along the x-axis, color-coded according to the classification of the arrays into the five cell cycle stages, M/G1 (yellow), G1 (green), S (blue), S/G2 (red), and G2/M (orange). The dashed unit and half-unit circles outline 100% and 25% of overall normalized array expression in the $|\alpha_1\rangle_N$ and $|\alpha_2\rangle_N$ subspace. (b) Correlation of each gene with $|\gamma_1\rangle_N$ vs. that with $|\gamma_2\rangle_N$, for 784 cell cycle regulated genes, color-coded according to the classification by Spellman *et al.* (3).

Alter, Orly *et al.* (2000) *Proc. Natl. Acad. Sci. USA* 97, 10101-10106