# Clustering Methods:
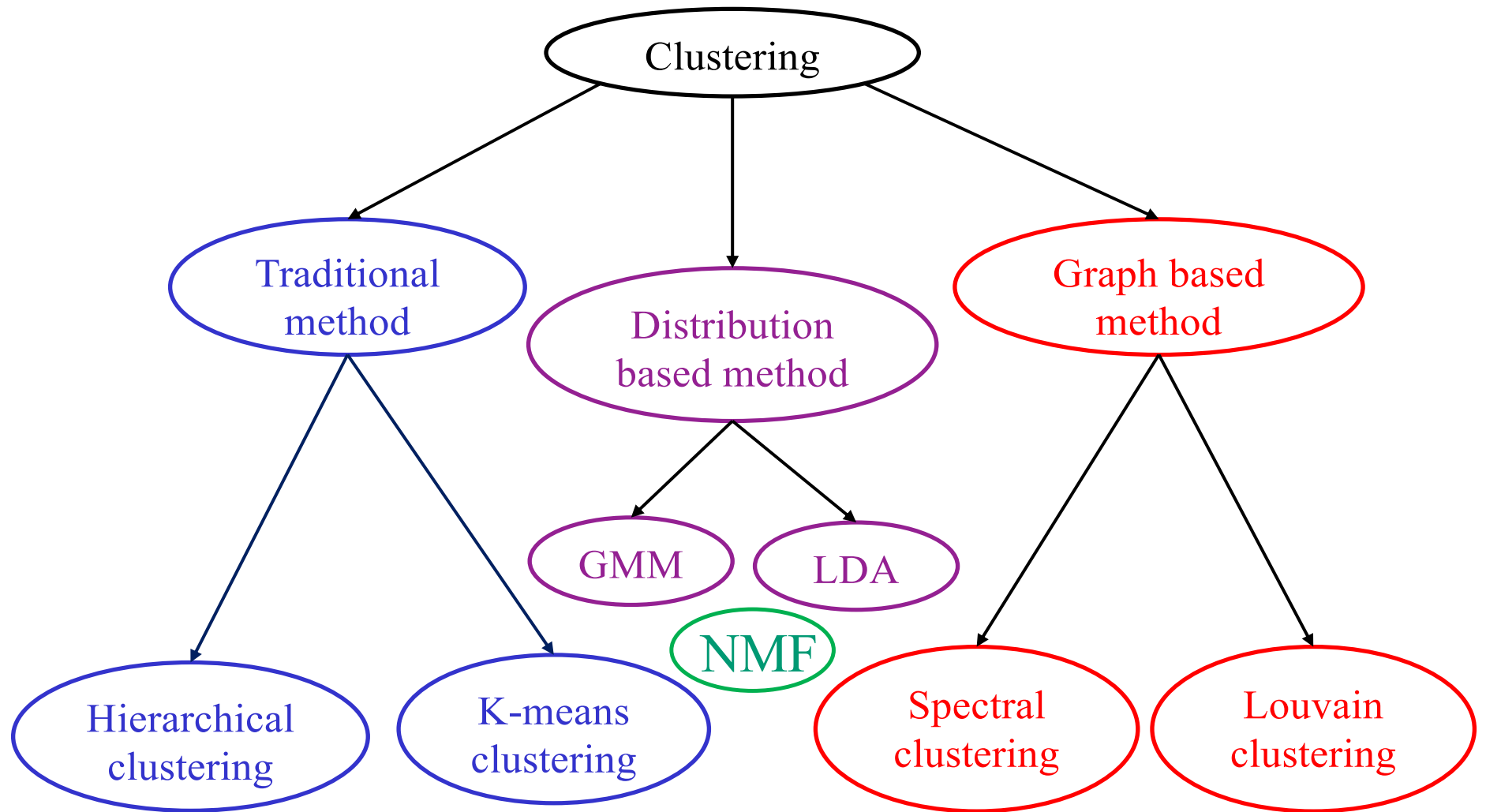# From k-means to Gaussian Mixture Model and Louvain Algorithm

Maxwell Lee

High-dimension Data Analysis Group
Laboratory of Cancer Biology and Genetics
Center for Cancer Research
National Cancer Institute

October 19, 2020
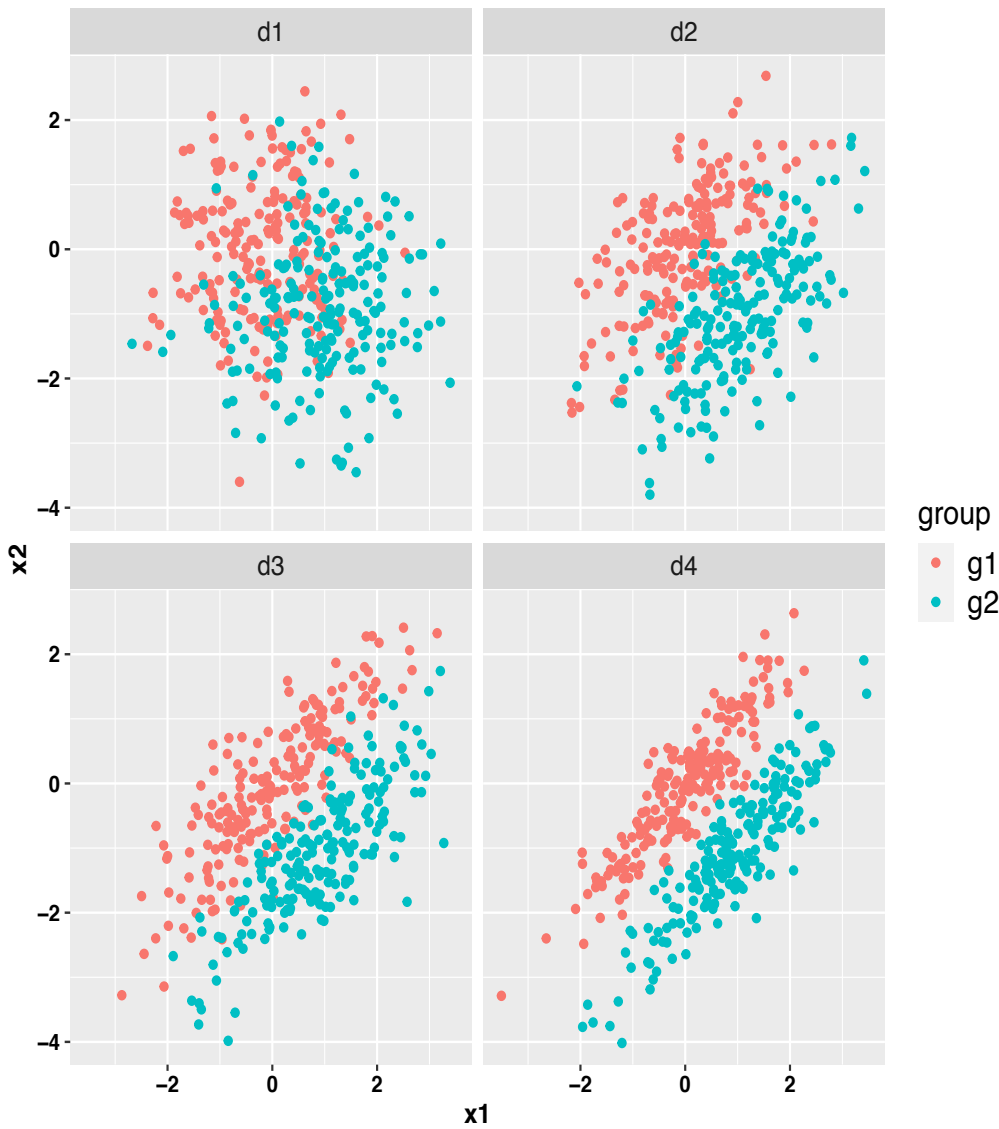
# Outline of Clustering Methods



GMM: Gaussian Mixture Model
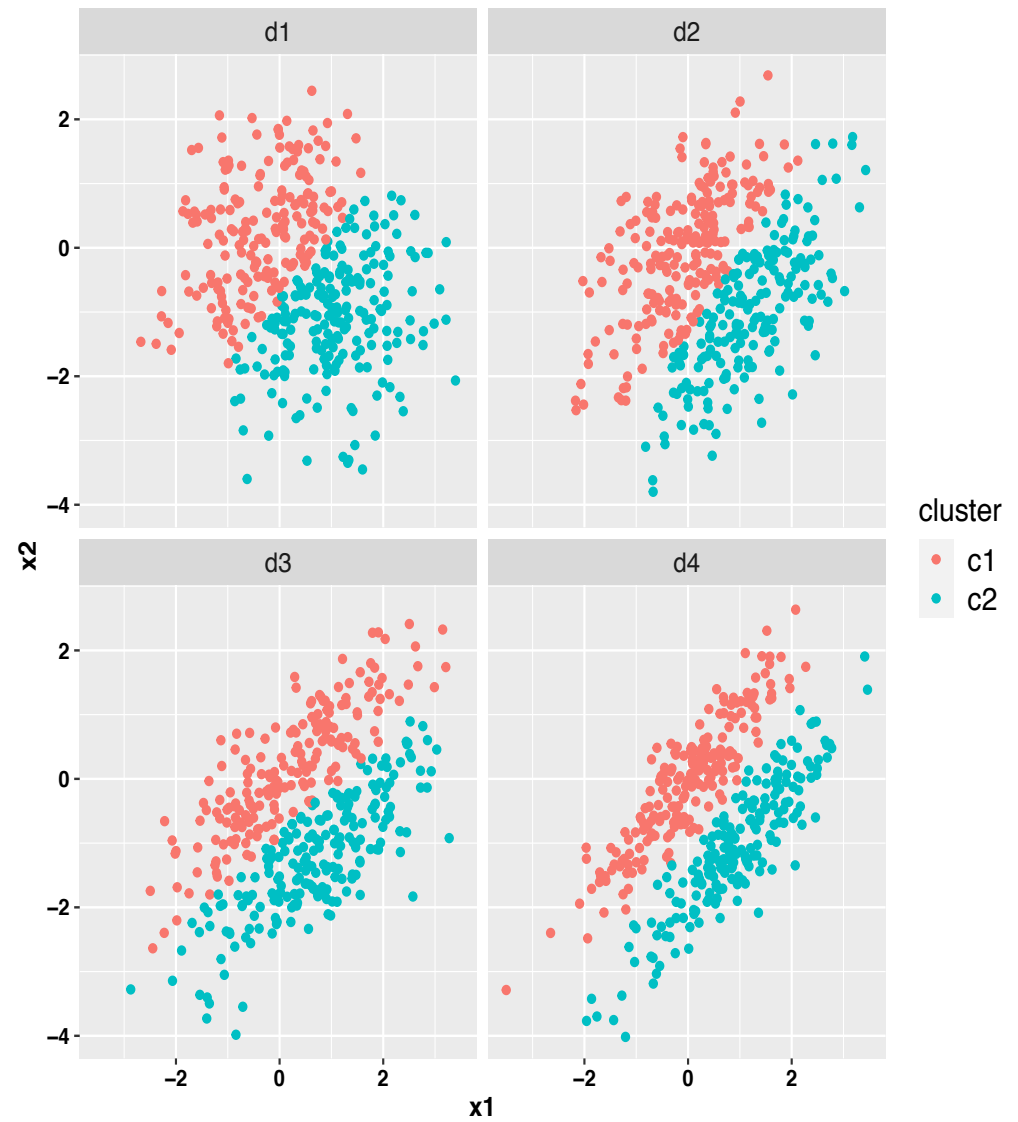LDA: Latent Dirichlet Allocation

NMF: Non-negative matrix factorization
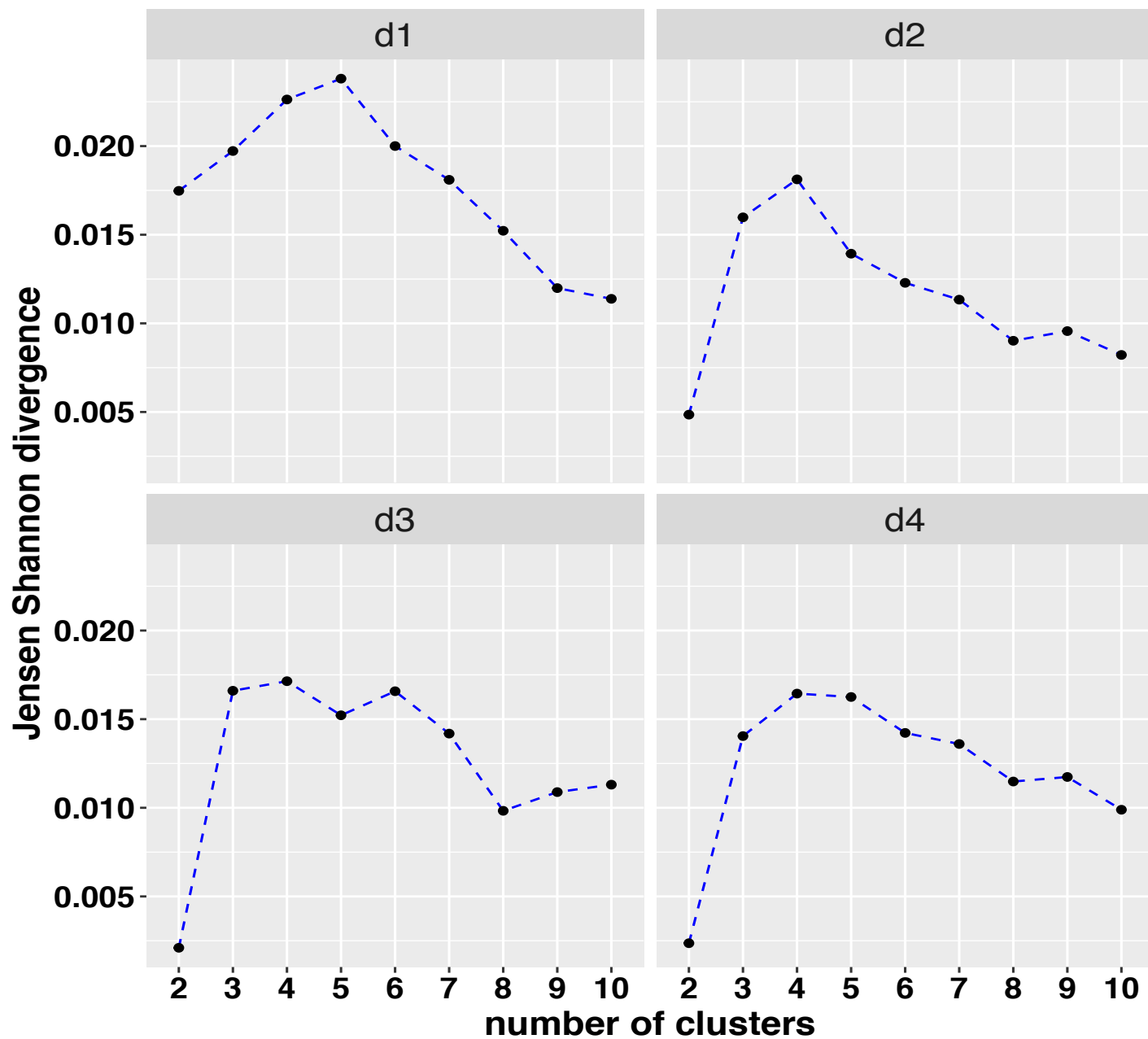
Contributed by Emily Tai

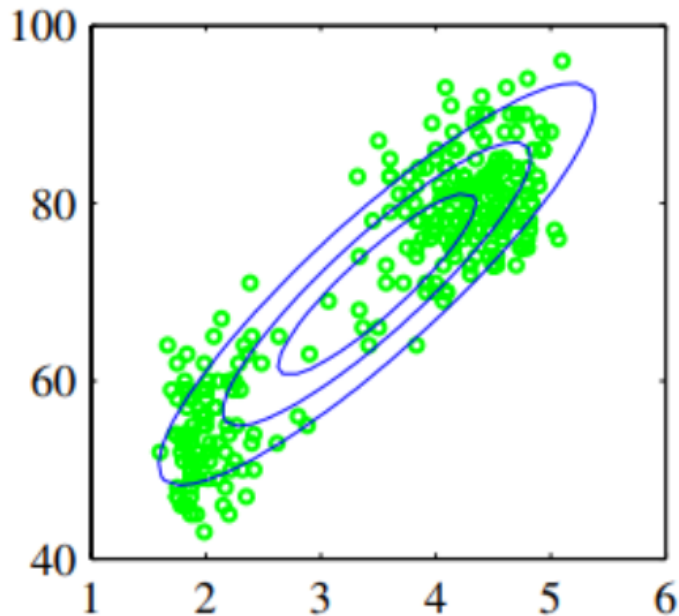Effect of Covariance Structure on GMM Clustering

# Choose the Number of Clusters with Jensen-Shannon Divergence
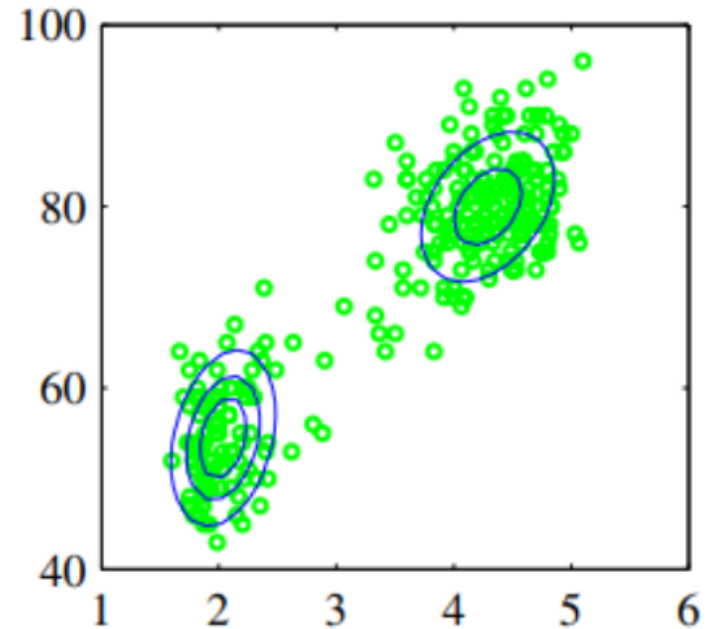
# Mixture of Bivariate Gaussian Distributions

## Single Gaussian



$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

## Mixture of two Gaussians



$$p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Component

Mixing coefficient

# Non-Negative Matrix Factorization (NMF)

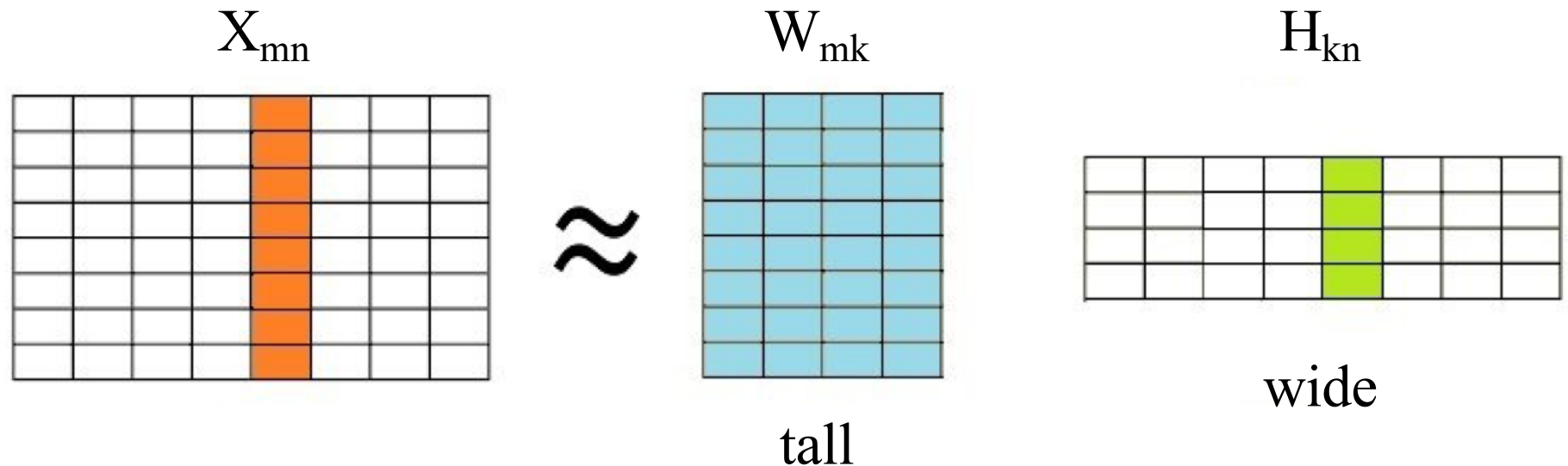NMF is a dimension reduction method. Why are we talking about it here?

Primary reason:
NMF is closely related to k-means and GMM. We can gain better understanding of clustering analysis through the lens of dimension reduction.

Secondary reason:
NMF is also closely related to Latent Dirichlet Allocation (LDA). It helps to understand LDA.

# Mathematic Model of Non-Negative Matrix Factorization

$X_{mn}$      $W_{mk}$      $H_{kn}$



$\approx$

tall

wide

$X_{mn}$: m features; n samples

$W_{mk}$: m features; k latent variables

$H_{km}$: k latent variables; n encodings

Each element of matrix is non-negative

X >= 0; W >= 0; H >= 0

latent variables:
basis images
topics
centroids
signatures

k << min(m,n) $\longrightarrow$ Dimension reduction

Lee and Seung, Nature 1999; 401:788–791

# NMF is Related to PCA

$X_{mn}$ $\approx$ $W_{mk}$ $H_{kn}$

$$X \sim WH$$

PCA

$$XE = Z$$
$$X = ZE^T$$
$$X^T = EZ^T$$
$$X^T \sim E_{mk}Z^T$$

But PCA can have negative values!

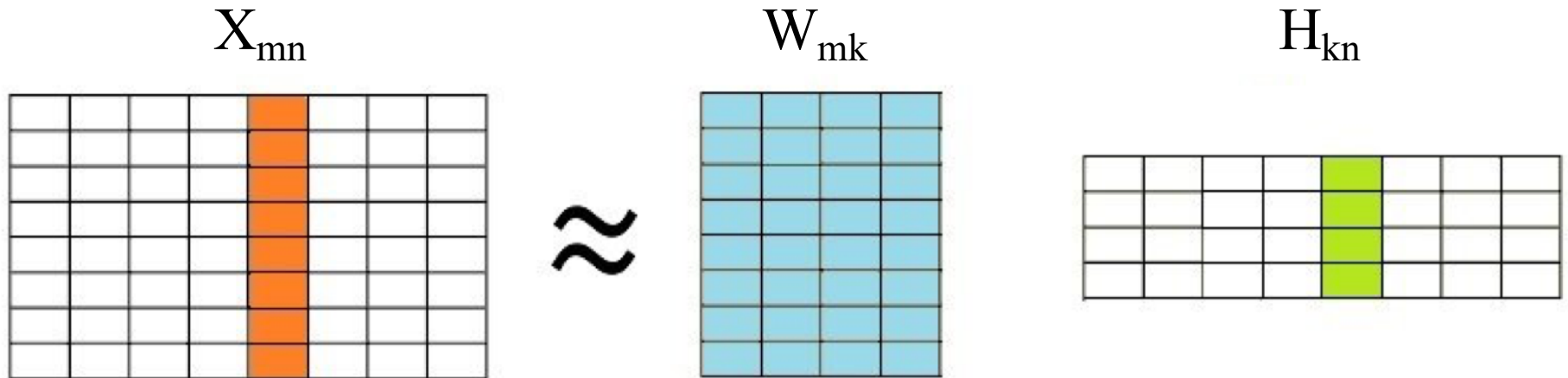# Understanding NMF from Topic Modeling (Mixture Model)

$$X_{mn} \qquad W_{mk} \qquad H_{kn}$$



$X_{mn}$: m words; n documents
$W_{mk}$: m words; k topics
$H_{km}$: k topics; n documents

Topic modeling:
$\Sigma_i X_{ij} = 1$ (column sums to 1)
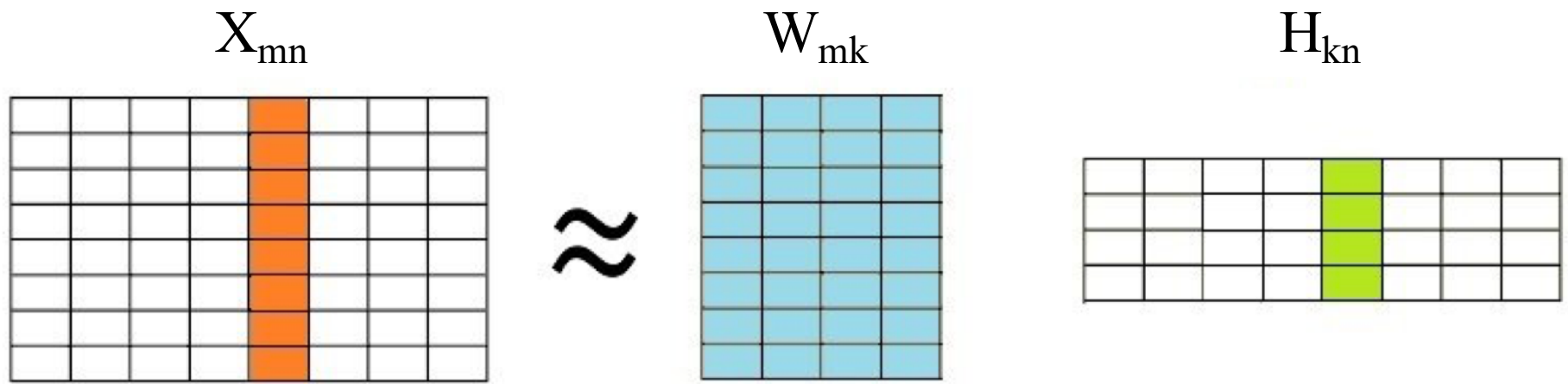$\Sigma_i W_{ij} = 1$ and $\Sigma_i H_{ij} = 1$

$x_{ij} \sim \Sigma_k w_{ik} h_{kj}$

# Understanding NMF from Topic Modeling (Mixture Model)

I will talk about spectral clustering, which is a graph-based method and consists of dimension reduction with Laplacian Eigenmap and k-means clustering in the reduced dimension space. I will also talk about Louvain algorithm, which is used in Seurat package to cluster single cell RNAseq data. Louvain algorithm is a network community approach. It is very fast and has capacity to do clustering analysis for million nodes in a network. I will provide practical examples to illustrate how each method works and how to interpret the results of clustering analysis and explain the pros and cons of each method.

# Understanding NMF from Topic Modeling (Mixture Model)

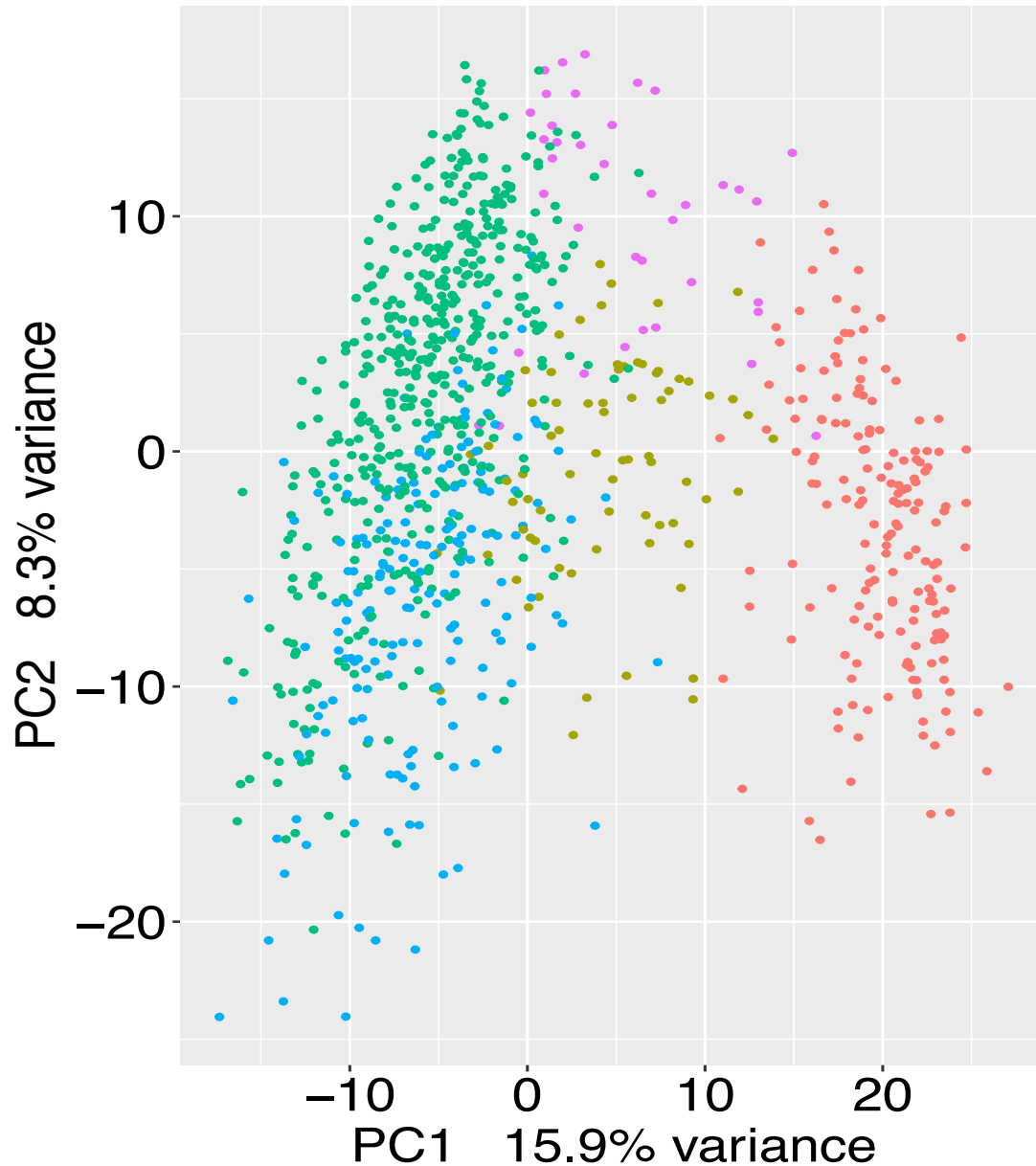$$X_{mn} \qquad W_{mk} \qquad H_{kn}$$



$$x_j \sim W\, h_j$$

Clustering:
label document with the topic having the highest frequency

nsNMF: nonsmooth NMF (sparse NMF)

PCA of TCGA BRCA Samples

$X_{mn}$
m = 5000 genes
n = 977 samples

PC2  8.3% variance
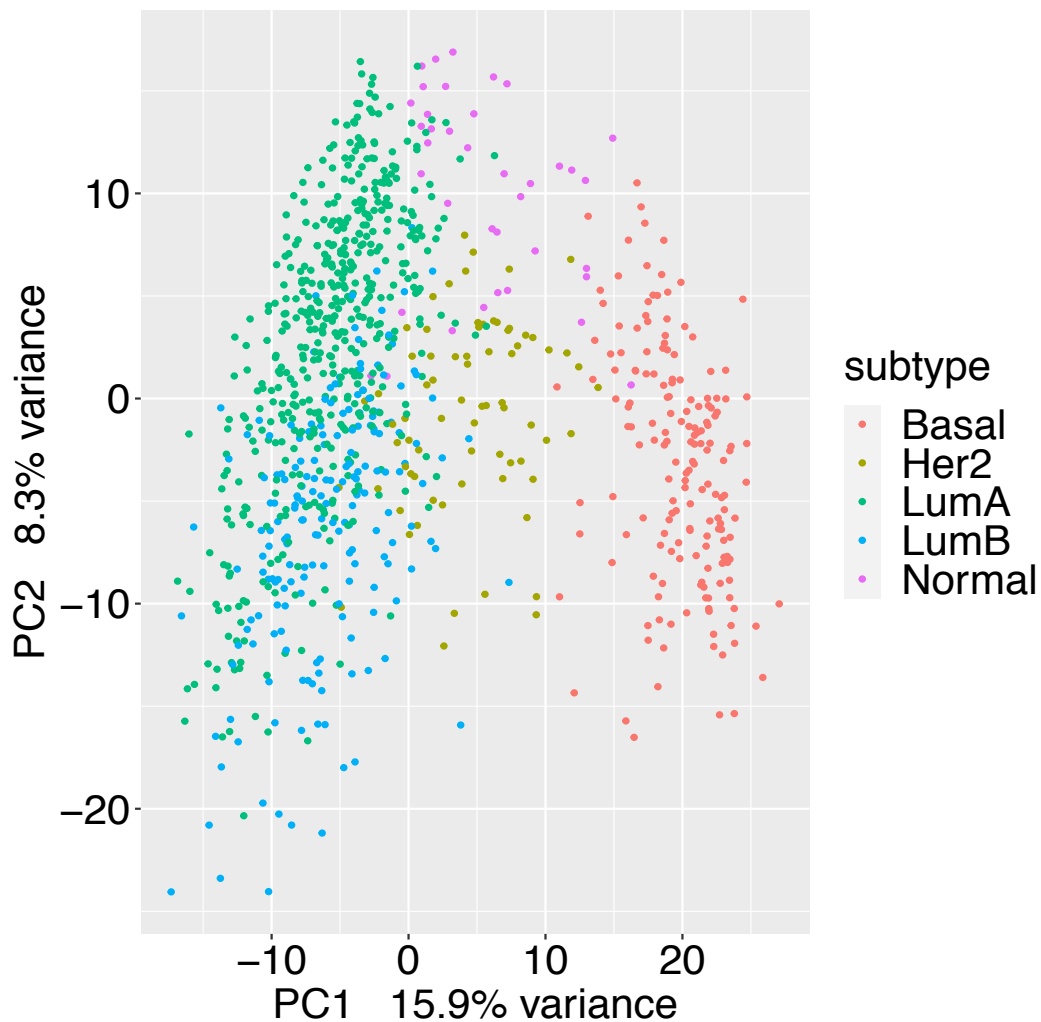
PC1    15.9% variance

subtype
- Basal
- Her2
- LumA
- LumB
- Normal

| x | freq |
|---|---|
| Basal | 173 |
| Her2 | 73 |
| LumA | 500 |
| LumB | 193 |
| Normal | 38 |

# PCA: Label by Subtype vs. by NMF Cluster



HD: high dimension, 5000 genes

# Comparison Between Subtype and NMF vs. k-means Cluster (HD)

|  | Basal | Her2 | LumA | LumB | Normal |
|---|---|---|---|---|---|
| clust1 | 155 | 2 | 11 | 1 | **30** |
| clust2 | 18 | 62 | 7 | 32 | 2 |
| clust3 | 0 | 4 | **338** | 38 | 6 |
| clust4 | 0 | 2 | 132 | 114 | 0 |
| clust5 | 0 | 3 | 12 | 8 | 0 |

NMF

Accuracy = (155 + 62 + 338 + 114) / 977 = 68.5%

<span style="color:red">Match</span>
<span style="color:green">Mismatch</span>

|  | Basal | Her2 | LumA | LumB | Normal |
|---|---|---|---|---|---|
| clust1 | 169 | 0 | 0 | 0 | 6 |
| clust2 | 4 | 69 | 17 | 40 | 5 |
| clust3 | 0 | 0 | 268 | 11 | 21 |
| clust4 | 0 | 0 | 125 | 119 | 0 |
| clust5 | 0 | 4 | **90** | 23 | 6 |

K-means

Accuracy = (169 + 69 + 268 + 119 + 6) / 977 = 64.6%

PCA of TCGA BRCA Samples with Pam50 Genes

977 samples
39 genes

subtype
Basal
Her2
LumA
LumB
Normal

PC2 21.1% variance

PC1 42.4% variance

PCA with pam50: Label by Subtype vs. by NMF Clusters

# PCA with pam50: Label by NMF vs. k-means Clusters



Label by NMF clusters in high-dimension

accuracy 69%

Label by k-means clusters in high-dimension
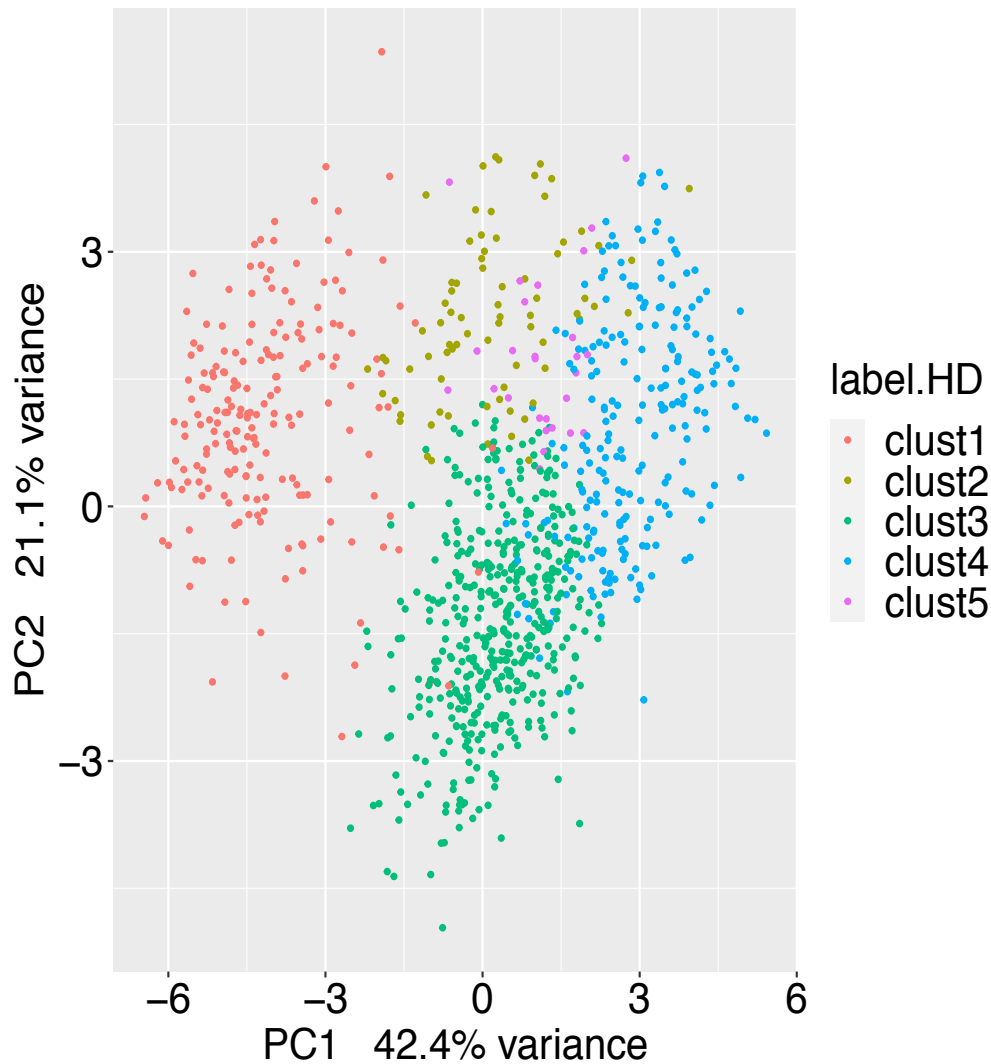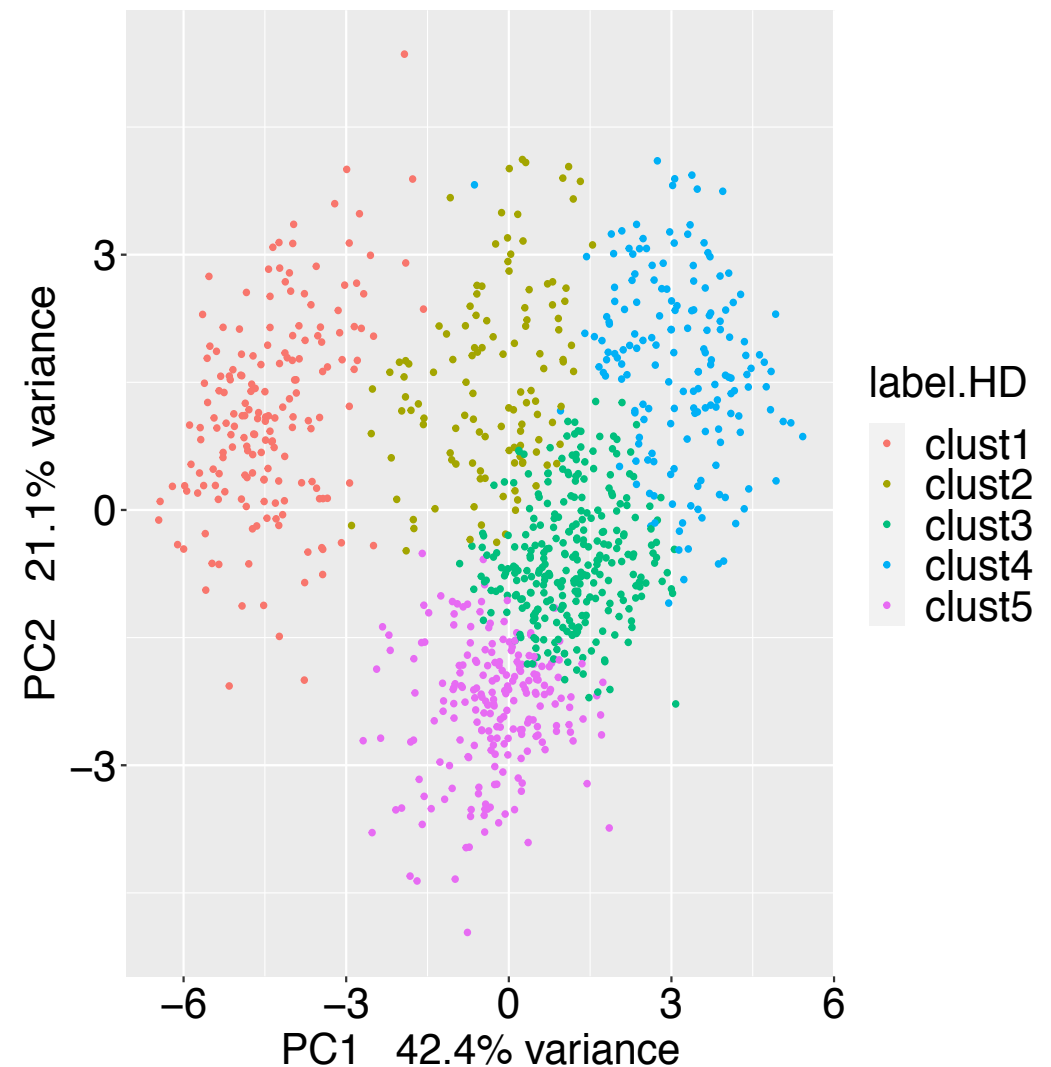
accuracy 59%

# Comparison Between Subtype and NMF vs. k-means Cluster (HD)

|        | Basal | Her2 | LumA | LumB | Normal |
|--------|-------|------|------|------|--------|
| clust1 | 172   | 6    | 5    | 0    | 15     |
| clust2 | 0     | 58   | 3    | 15   | 0      |
| clust3 | 0     | 9    | **356** | 61 | 23   |
| clust4 | 0     | 0    | **136** | 91 | 0    |
| clust5 | 1     | 0    | 0    | **26** | 0    |

NMF

Accuracy = (172 + 58 + 356+ 91) / 977 = 69%

Match
Mismatch

|        | Basal | Her2 | LumA | LumB | Normal |
|--------|-------|------|------|------|--------|
| clust1 | 170   | 0    | 0    | 0    | 8      |
| clust2 | 2     | 72   | 11   | 31   | 4      |
| clust3 | 0     | 0    | 221  | 75   | 0      |
| clust4 | 1     | 1    | 81   | 87   | 0      |
| clust5 | 0     | 0    | **187** | 0  | **26** |

K-means

Accuracy = (170 + 72 + 221 + 87 + 26) / 977 = 59%

# Application of NMF to COSMIC Mutation Signature

|   | A | C | G | T |
|---|---|---|---|---|
| A | 0 | *1* | 1 | 1 |
| C | 1 | 0 | 1 | 1 |
| G | 1 | 1 | 0 | 1 |
| T | 1 | 1 | 1 | 0 |

$N_1[A>C]_2N_3$

$4*6*4=96$

# COSMIC Mutation Signatures



NMF analysis of 7042 tumors                    15 out of 22 signatures

Ludmil B. Alexandrov, …, Michael R. Stratton Nature 500:415–421(2013)

# COSMIC Mutation Signatures

# COSMIC Mutation Signatures



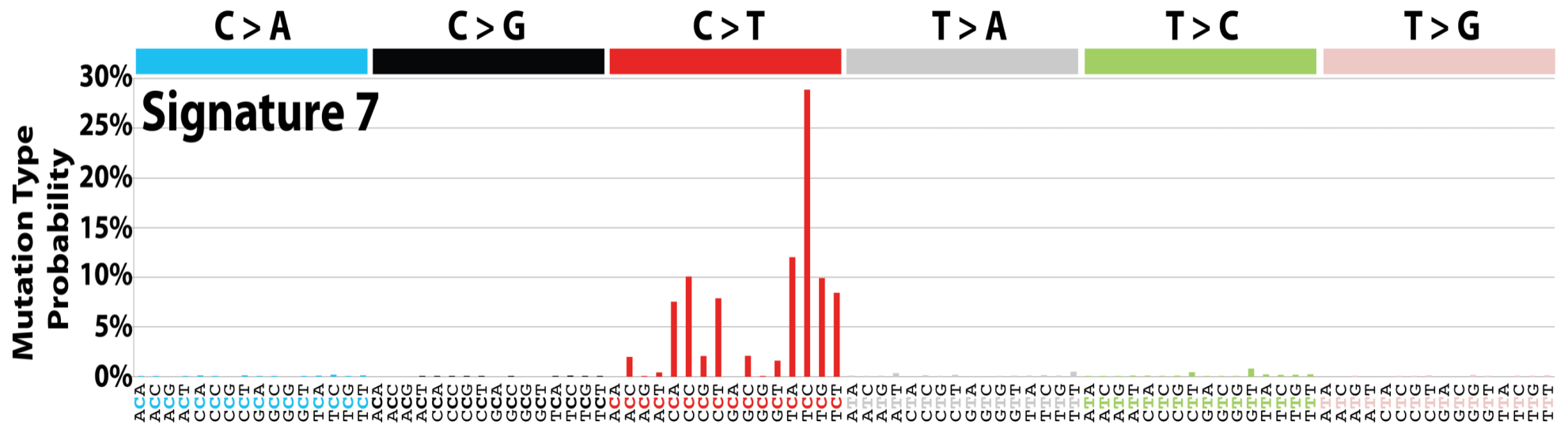Ludmil B. Alexandrov, …, Michael R. Stratton Nature 500:415–421(2013)

# Algorithm of NMF

$$X \approx WH$$

$$\min_{W,H \geq 0} ||X - WH||_F^2 = \sum_{i,j}(X - WH)_{ij}^2$$

Alternating multiplicative update

$$H \leftarrow H \odot W^T X / W^T WH$$

$$W \leftarrow W \odot X H^T / WHH^T$$

Hadamard product operator $\odot$
element-wise multiplication of matrices

# Matrix Representation of k-means Clustering

$$\min_{C_1,C_2,\ldots,C_k} \sum_{i=1}^{k} \sum_{\mathbf{x}\in C_i} ||\mathbf{x} - \mu_i||^2$$

$\mu_i$ is the centroid of data points in $C_i$

$$B_{nk} = \begin{array}{c|cccc} & C1 & C2 & C3 & C4 \\ \hline x_1 & 1 & 0 & 0 & 0 \\ x_2 & 0 & 1 & 0 & 0 \\ x_3 & 0 & 1 & 0 & 0 \\ x_4 & 0 & 0 & 1 & 0 \\ x_5 & 0 & 0 & 0 & 1 \\ x_6 & 0 & 0 & 1 & 0 \end{array}$$

Class membership matrix B:
each row has only one 1; the others are 0
sum of column is the size of the cluster
columns are orthogonal

# Matrix Representation of k-means Clustering

$B_{nk} =$

|       | C1 | C2 | C3 | C4 |
|-------|----|----|----|----|
| $x_1$ | 1  | 0  | 0  | 0  |
| $x_2$ | 0  | 1  | 0  | 0  |
| $x_3$ | 0  | 1  | 0  | 0  |
| $x_4$ | 0  | 0  | 1  | 0  |
| $x_5$ | 0  | 0  | 0  | 1  |
| $x_6$ | 0  | 0  | 1  | 0  |

$\sum_i B_{ij} = |\mathcal{C}_j|$   sum of column is the size of the cluster

$$D = diag(1/|\mathcal{C}_1|, 1/|\mathcal{C}_2|, ..., 1/|\mathcal{C}_k|)$$

$$(BD^{\frac{1}{2}})^T BD^{\frac{1}{2}} = D^{\frac{1}{2}} B^T BD^{\frac{1}{2}} = I$$

# Matrix Representation of k-means Clustering

$$\min_{\mathcal{C}_1,\mathcal{C}_2,\ldots,\mathcal{C}_k} \sum_{i=1}^{k} \sum_{\mathbf{x}\in\mathcal{C}_i} ||\mathbf{x} - \mu_i||^2$$

$$\mu_i = \frac{1}{|\mathcal{C}_i|} \sum_{\mathbf{x}\in\mathcal{C}_i} \mathbf{x}$$

M = XBD

M: k means (centroids) in columns
X: n samples in columns
B: k clusters in columns
D: diagonal matrix; 1/cluster size

$XBDB^T$: each sample selects its corresponding centroid in columns

$$\min_{B} ||X - XBDB^T||_F^2$$

# k-means Clustering is Equivalent to Sparse NMF

k-means
$$\min_{B} ||X - XBDB^T||_F^2$$

NMF
$$\min_{W,H \geq 0} ||X - WH||_F^2$$

W ~ XBD
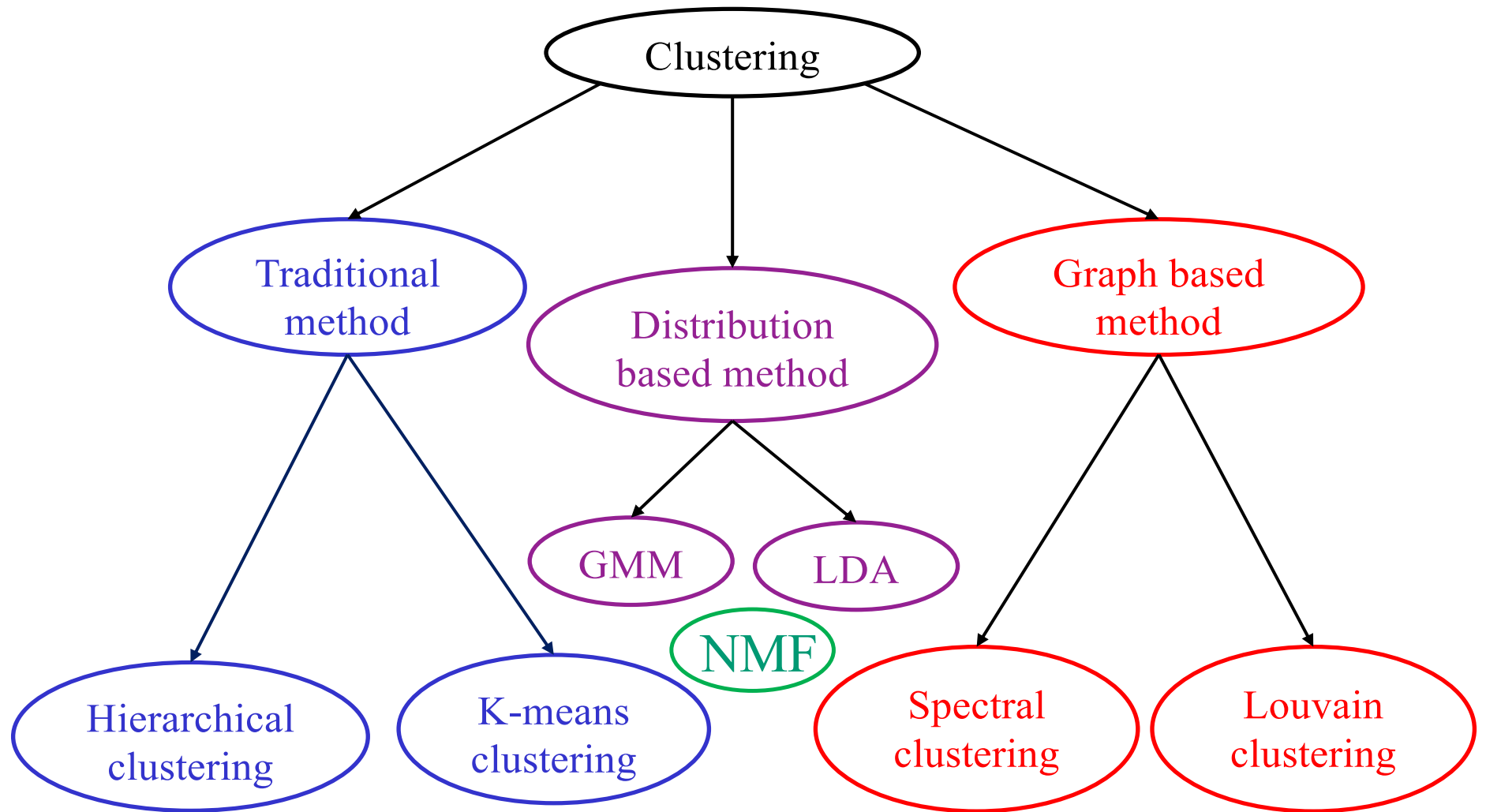
H ~ B$^T$

B is orthogonal and sparse

Sparse NMF or non-smooth NMF

NMF is related to GMM and mixture model

# Outline of Clustering Methods



GMM: Gaussian Mixture Model
LDA: Latent Dirichlet Allocation

NMF: Non-negative matrix factorization

Contributed by Emily Tai