

Machine learning ideas in Seurat packages

Jan/16/2020

Sunho Park

Outline

- **Dimensionality reduction**

- PCA
- t-SNE
- Uniform Manifold Approximation and Projection (UMAP)

- **Clustering approaches**

- Graph based methods: e.g., Louvain algorithm
- Cell-type identification

- **Multiple Dataset Integration**

- Canonical correlation analysis & L2-norm normalization
- Anchoring

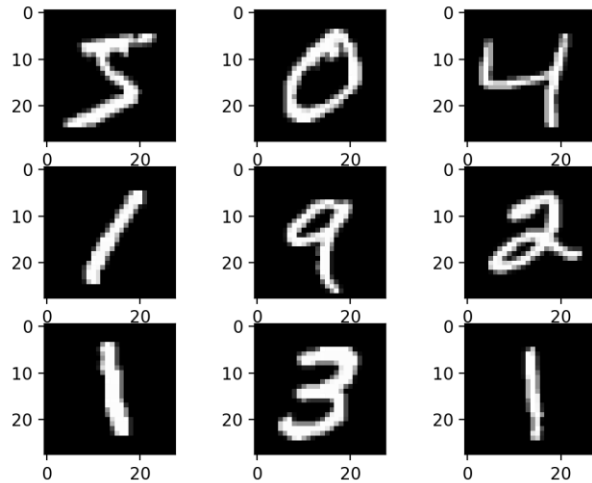
Outline

- **Dimensionality reduction (visualization)**
 - PCA
 - t-SNE
 - Uniform Manifold Approximation and Projection (UMAP)
- Clustering approaches
 - Graph based methods: Louvain algorithm
 - Cell-type identification
- Multiple Dataset Integration
 - Canonical correlation analysis & L2-norm normalization
 - Anchoring

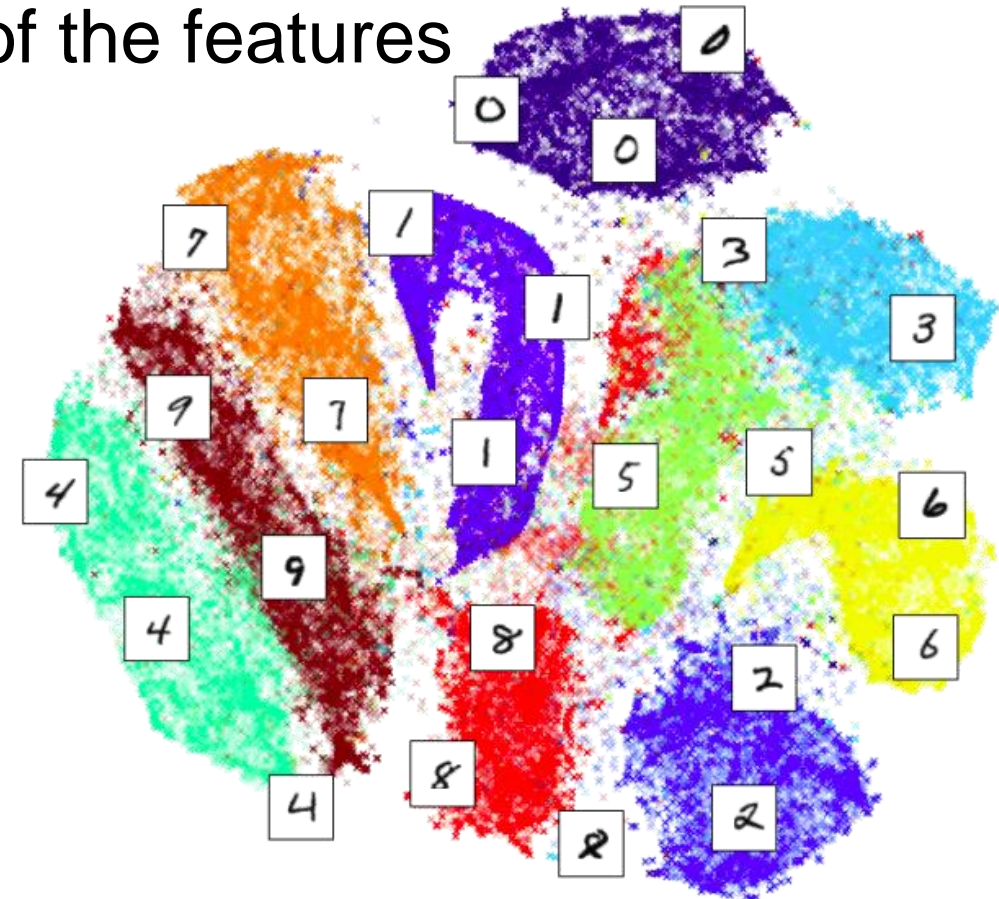
Dimensionality reduction

Figure credit: liam schoneveld

- A process of reducing the number of the features



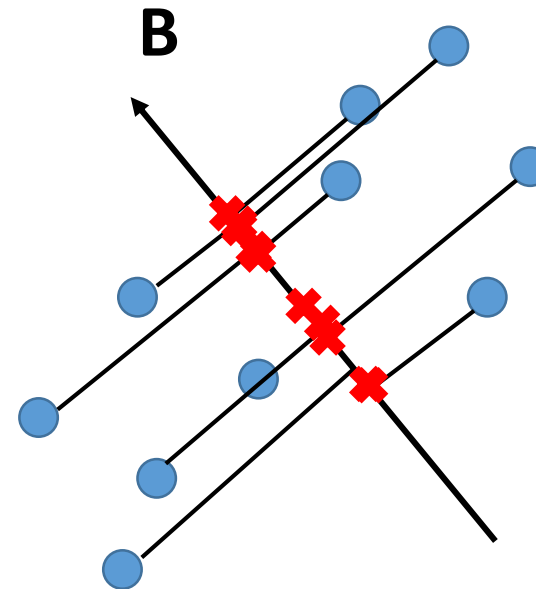
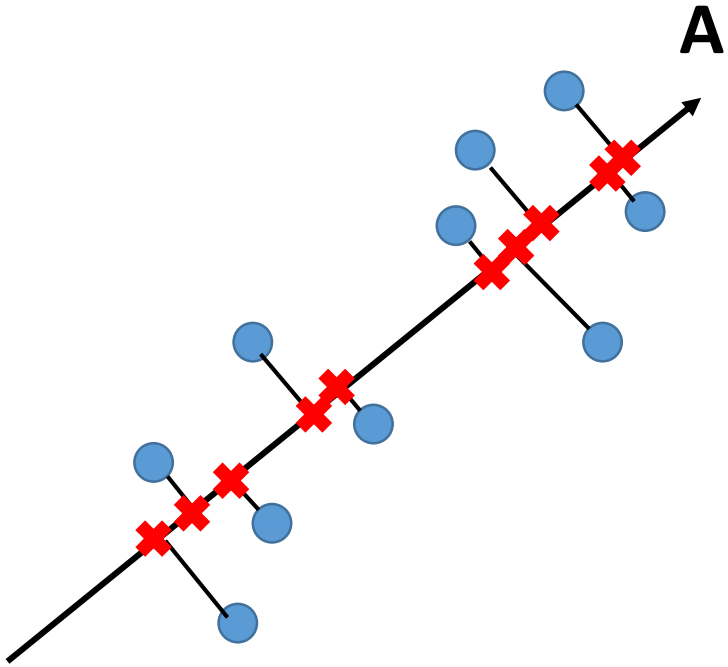
Original image (MNIST): $\mathbf{x}_i \in \mathbb{R}^{784}$ (28×28)



Reduced space (by t-SNE): $\mathbf{x}_i \in \mathbb{R}^2$

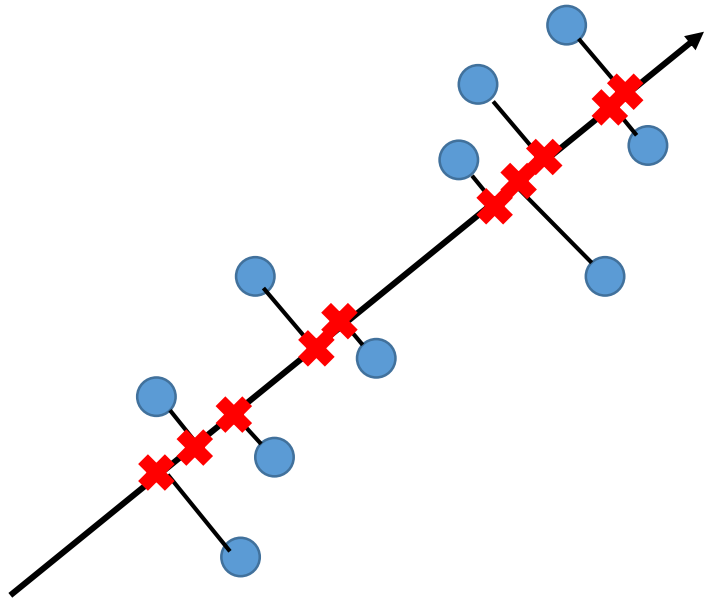
Dimensionality reduction: PCA

- Which one is a more accurate lower dimensional representation?



Dimensionality reduction: PCA

- Find a direction that maximizes the projected variance



- Data samples are assumed to be centered

$$\tilde{\mathbf{x}}_i = \mathbf{x}_i - \bar{\mathbf{x}}$$

- Maximizing the projected variance

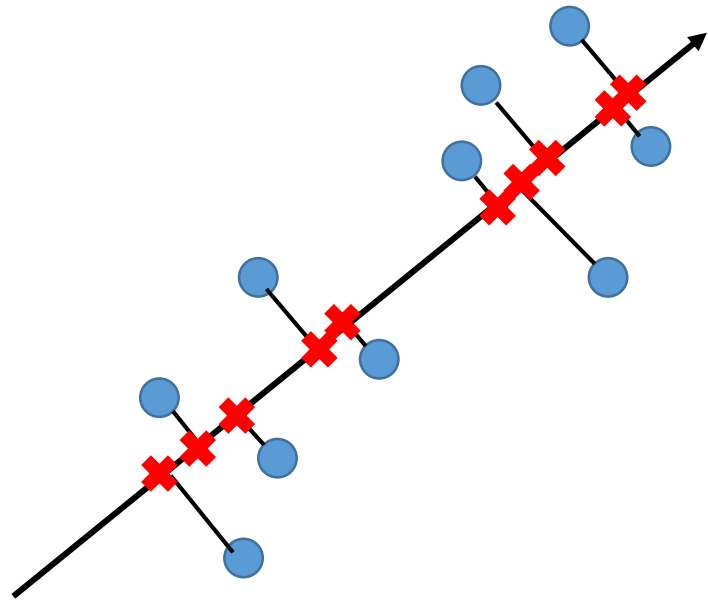
$$\hat{\mathbf{u}} = \operatorname{argmax}_{\mathbf{u}} \frac{1}{N} \sum_{i=1}^N (\mathbf{u}^T \tilde{\mathbf{x}}_i)^2$$

$$= \operatorname{argmax}_{\mathbf{u}} \mathbf{u}^T \left(\frac{1}{N} \sum_{i=1}^N \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T \right) \mathbf{u}$$

$$\text{such that } \|\mathbf{u}\|_2^2 = 1$$

Dimensionality reduction: FLD

- Find a direction that maximizes the projected variance



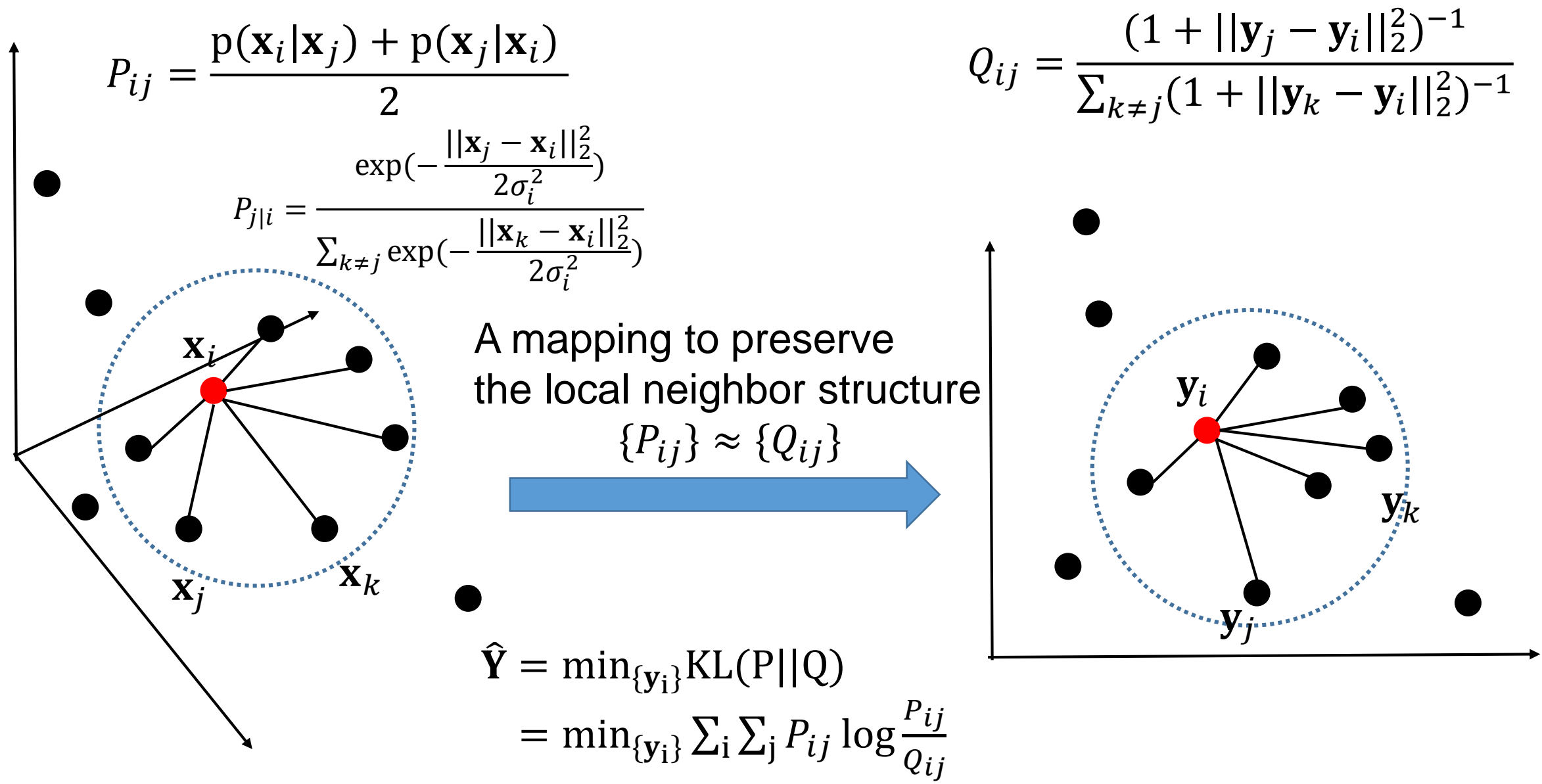
$$\hat{\mathbf{u}} = \operatorname{argmax}_{\mathbf{u}} \mathbf{u}^T \mathbf{S} \mathbf{u} \quad \mathbf{S} = \left(\frac{1}{N} \sum_{i=1}^N \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T \right)$$

$$\text{such that } \|\mathbf{u}\|_2^2 = 1$$

- $\hat{\mathbf{u}}$ is the first eigenvector (with the largest eigenvalue) of the sample covariance matrix \mathbf{S}
- If we need a d dimensional subspace,
 $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d]$ with $\lambda_1 \geq \lambda_2 \dots \geq \lambda_d$

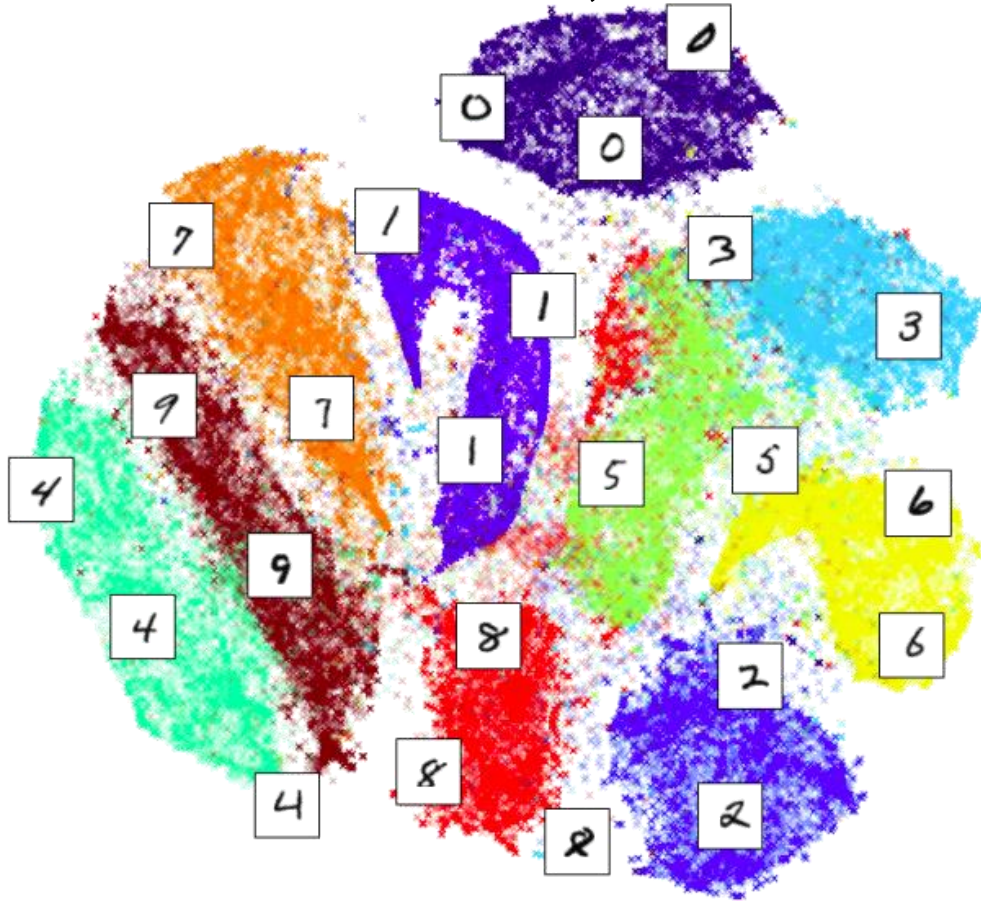
$$\mathbf{y}_i = \mathbf{U}^T \tilde{\mathbf{x}}_i \quad \mathbf{u}_i \perp \mathbf{u}_j \text{ for } i \neq j$$

Dimensionality reduction: t-SNE



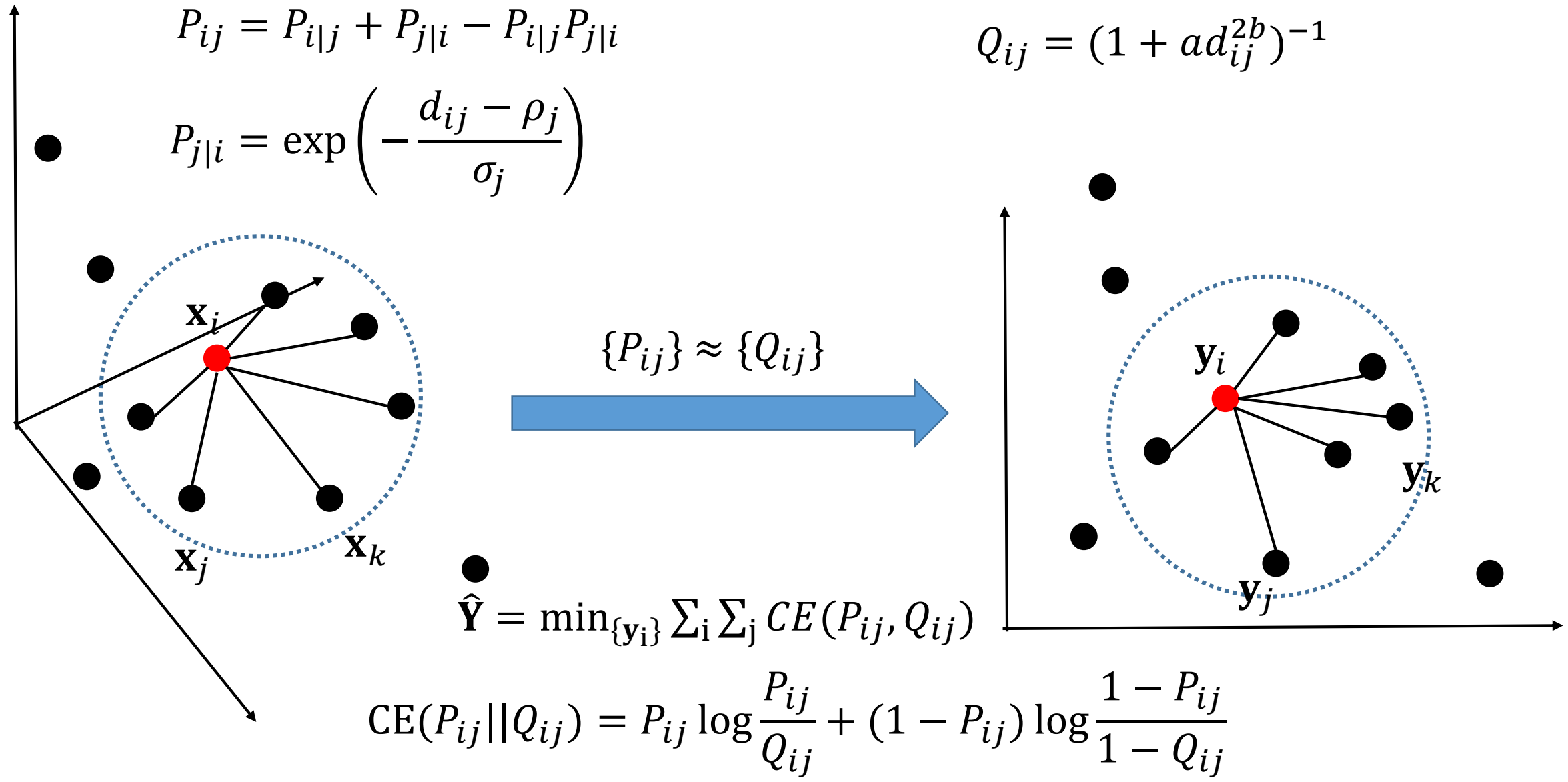
Dimensionality reduction

- t-SNE has been widely used for visualization of high dimensional data, but it has some limitations:



- **The global structure is not preserved**
 - points in the same cluster are close to each other (local structure)
 - Similarities between clusters might not be accurate (global structure, e.g., consider clusters (0 and 3), (0 and 6))
- **It is not scalable to handle fast growing sample sizes in single cell data.**

Uniform Manifold Approximation and Projection (UMAP)



Dimensionality reduction (t-SNE vs UMAP)

- UMAP works better for preserving the global structure of the data than t-SNE

- t-SNE

$$\text{KL}(P_{ij}||Q_{ij}) \approx \exp(-d_{ij}^X) \log(1 + (d_{ij}^Y)^2)$$

- d_{ij}^X is large: d_{ij}^Y can be any value (i.e., the global structure is not guaranteed)

- UMAP

$$\text{CE}(P_{ij}||Q_{ij}) \approx \exp(-d_{ij}^X) \log(1 + (d_{ij}^Y)^2) + (1 - \exp(-d_{ij}^X)) \log\left(\frac{1 + (d_{ij}^Y)^2}{(d_{ij}^Y)^2}\right)$$

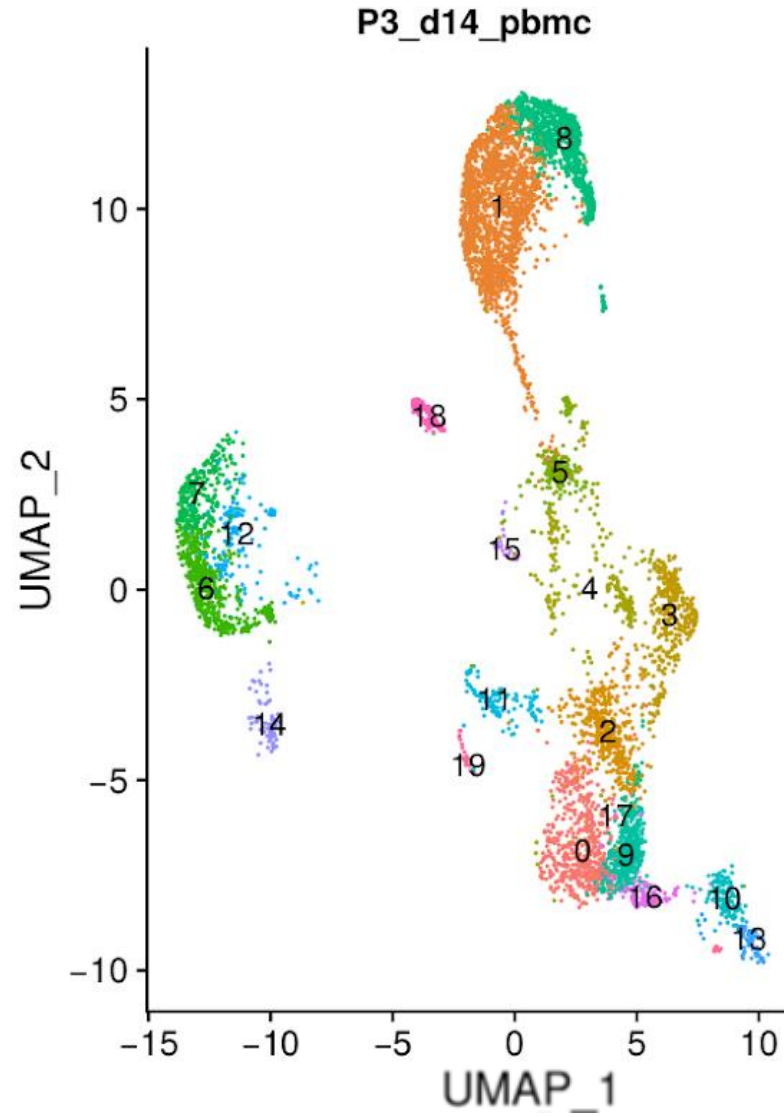
- d_{ij}^X is small: $\text{CE}(P_{ij}||Q_{ij}) \approx \log(1 + (d_{ij}^Y)^2)$

- d_{ij}^X is large: $\text{CE}(P_{ij}||Q_{ij}) \approx \log\left(\frac{1 + (d_{ij}^Y)^2}{(d_{ij}^Y)^2}\right)$ (i.e., it gives a high penalty for a small d_{ij}^Y and thus the global structure can be preserved)

Outline

- Dimensionality reduction (visualization)
 - PCA
 - t-SNE
 - Uniform Manifold Approximation and Projection (UMAP)
- **Clustering approaches**
 - Graph based methods: e.g., Louvain algorithm
 - Cell-type identification
- Multiple Dataset Integration
 - Canonical correlation analysis & L2-norm normalization
 - Anchoring

Clustering approaches



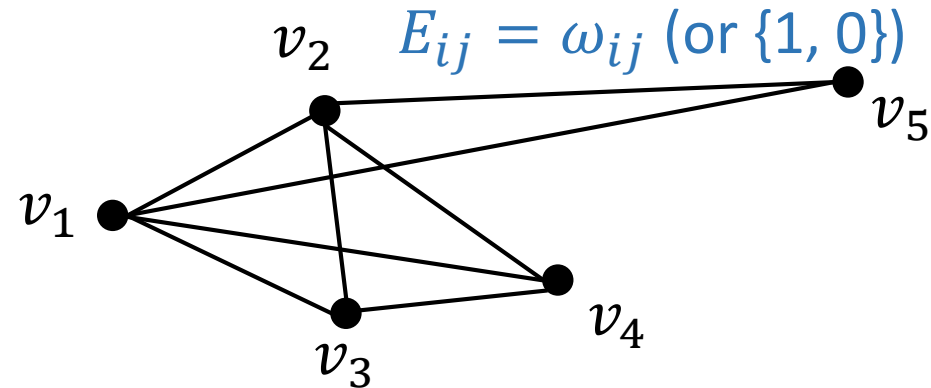
Community detection (clustering) in networks (graphs)

■ Graph

It is used to represent complex systems (e.g., friend networks on Facebook, gene-gene interaction networks)

- $V = \{v_1, v_2, \dots, v_N\}$: nodes (vertices)
- $E = \{e_{ij}\}$: links (edges)

$$G = \{V, E\}$$



Clustering on graphs

- Goal: divide a graph into multiple clusters where the nodes in the same cluster are more close to each other than to those in other clusters

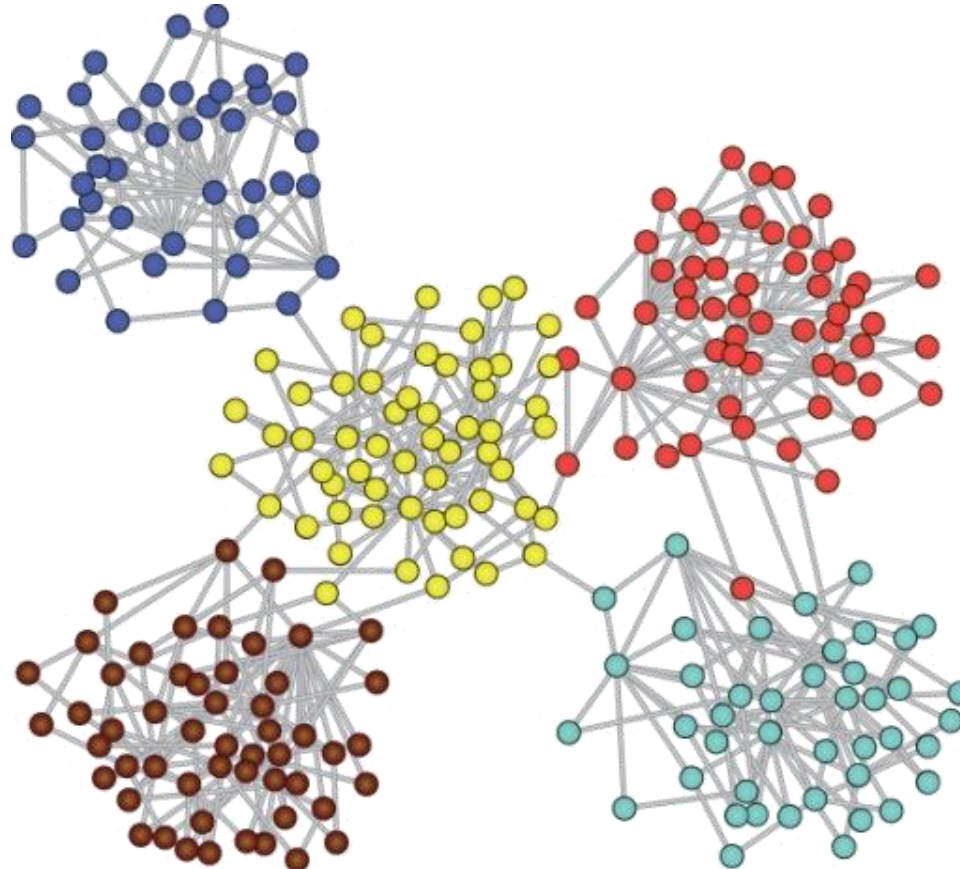


Figure credit: <https://github.com/benedekrozemberczki/awesome-community-detection>

Graph clustering

■ Modularity

- Measurement of the strength of partitioning nodes into modules (clusters)

$$Q(\gamma) = \frac{1}{2m} \sum_{i,j} \left[E_{ij} - \frac{d_i d_j}{2m} \right] \delta(\gamma(v_i), \gamma(v_j))$$

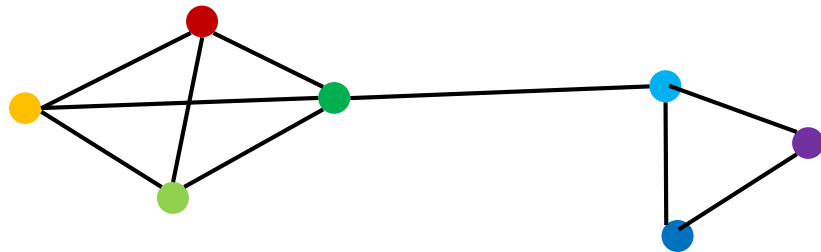
- d_i : degree of the node i ($d_i = \sum_j E_{ij}$) and $m = \frac{1}{2} \sum_{i,j} E_{ij}$
 - $\frac{d_i d_j}{2m}$: probability of an edge existing between v_i and v_j in the random null model
 - the fraction of the edges inside the cluster minus the expected fraction if edges were distributed at random
-
- The high modularity means that there are dense connections between the nodes in the same cluster but sparse connections between nodes in different clusters

Graph clustering

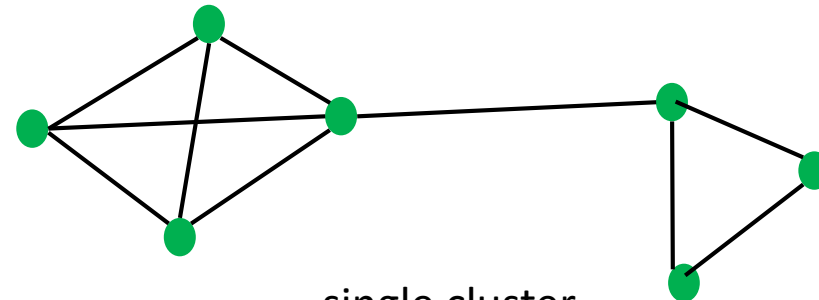
■ Modularity

$$Q(\gamma) = \frac{1}{2m} \sum_{i,j} \left[E_{ij} - \frac{d_i d_j}{2m} \right] \delta(\gamma(v_i), \gamma(v_j))$$

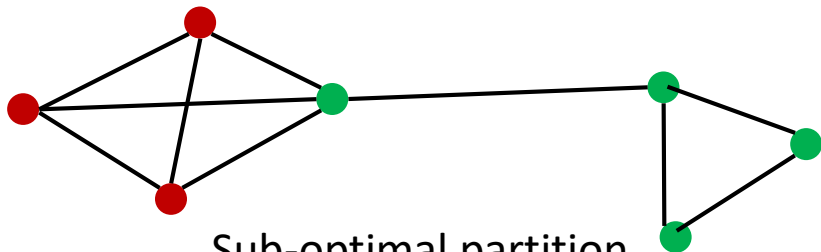
The high modularity means that there are dense connections between the nodes in the same cluster but sparse connections between nodes in different clusters



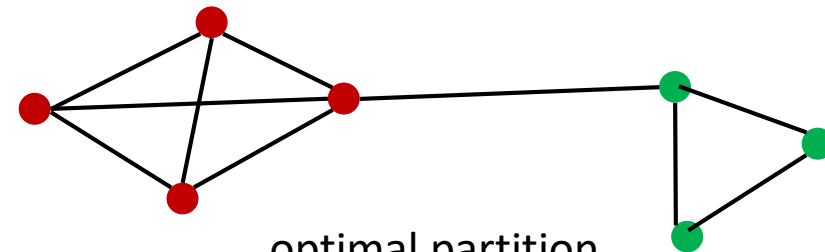
of clusters = # of nodes



single cluster



Sub-optimal partition



optimal partition

Clustering on graphs

- Louvain algorithm
 - greedily maximizes the modularity score by using an agglomerative approach
 - is one of the fastest modularity-based algorithms and scalable

Fast unfolding of communities in large networks

Vincent D. Blondel^{1;a}, Jean-Loup Guillaume^{1,2;b}, Renaud Lambiotte^{1,3;c} and Etienne Lefebvre¹

¹Department of Mathematical Engineering, Université catholique de Louvain, 4 avenue Georges Lemaitre, B-1348 Louvain-la-Neuve, Belgium

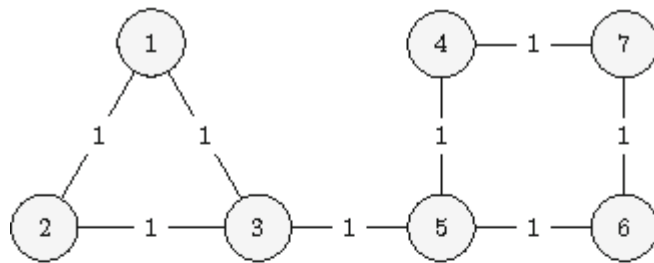
² LIP6, Université Pierre et Marie Curie, 4 place Jussieu, 75005 Paris, France

³ Institute for Mathematical Sciences, Imperial College London, 53 Prince's Gate, South Kensington campus, SW72PG, UK

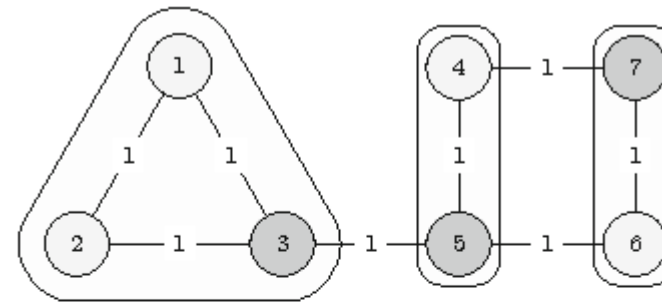
E-mail: ^avincent.blondel@uclouvain.be; ^bjean-loup.guillaume@lip6.fr;

^cr.lambiotte@imperial.ac.uk;

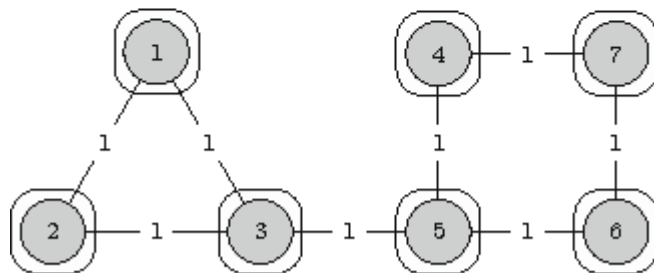
Clustering on graphs (Louvain algorithm)



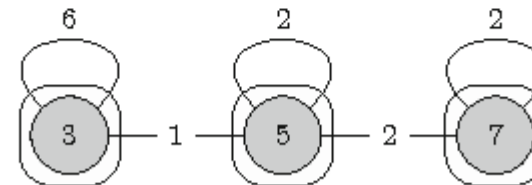
(a) original network



(c) step 1 of 1st iteration



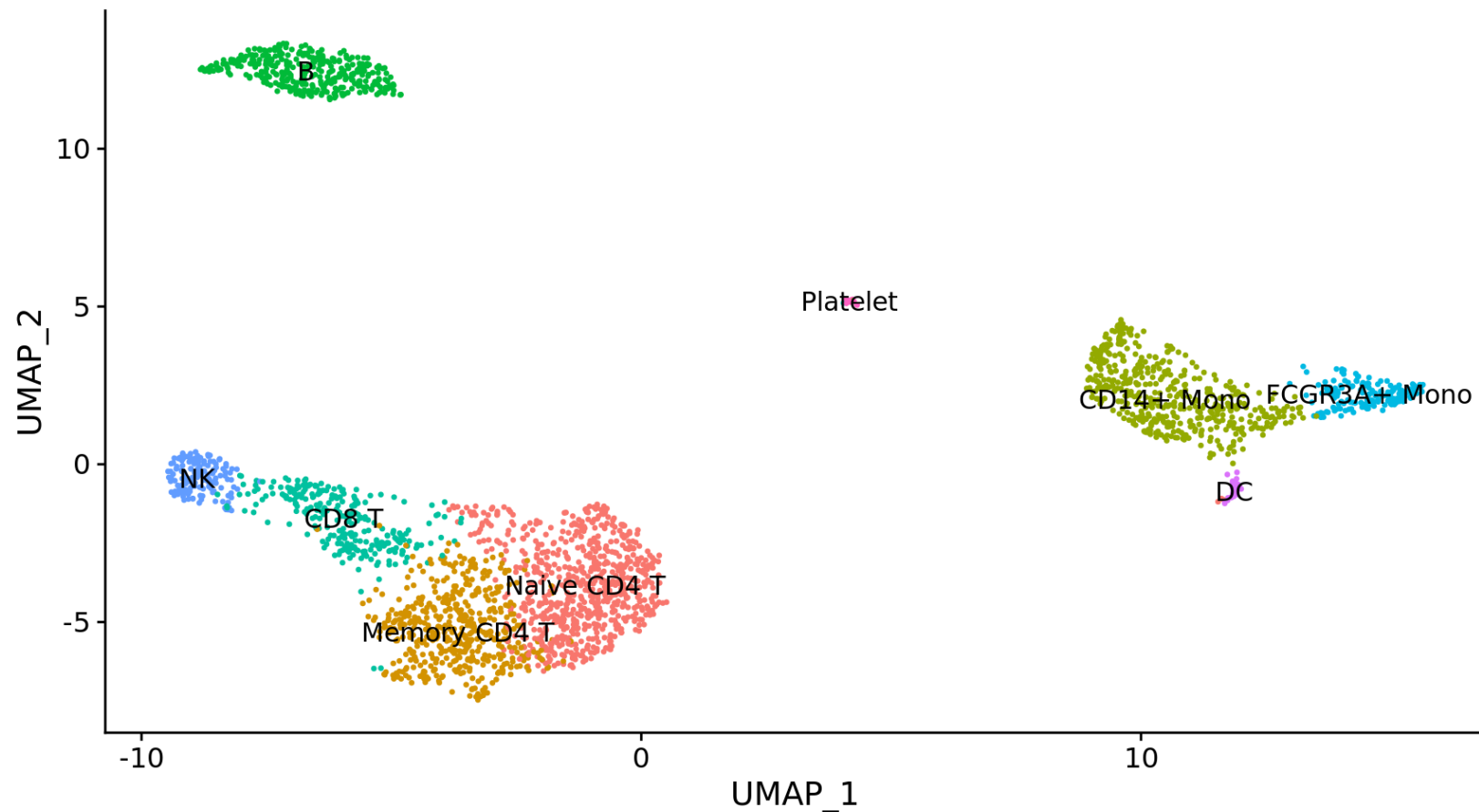
(b) initial communities



(d) step 2 of 1st iteration

Cell-type identification

- Identify cell-types from single-cell RNA sequencing data



Cell-type identification

▪ **Clustering step**

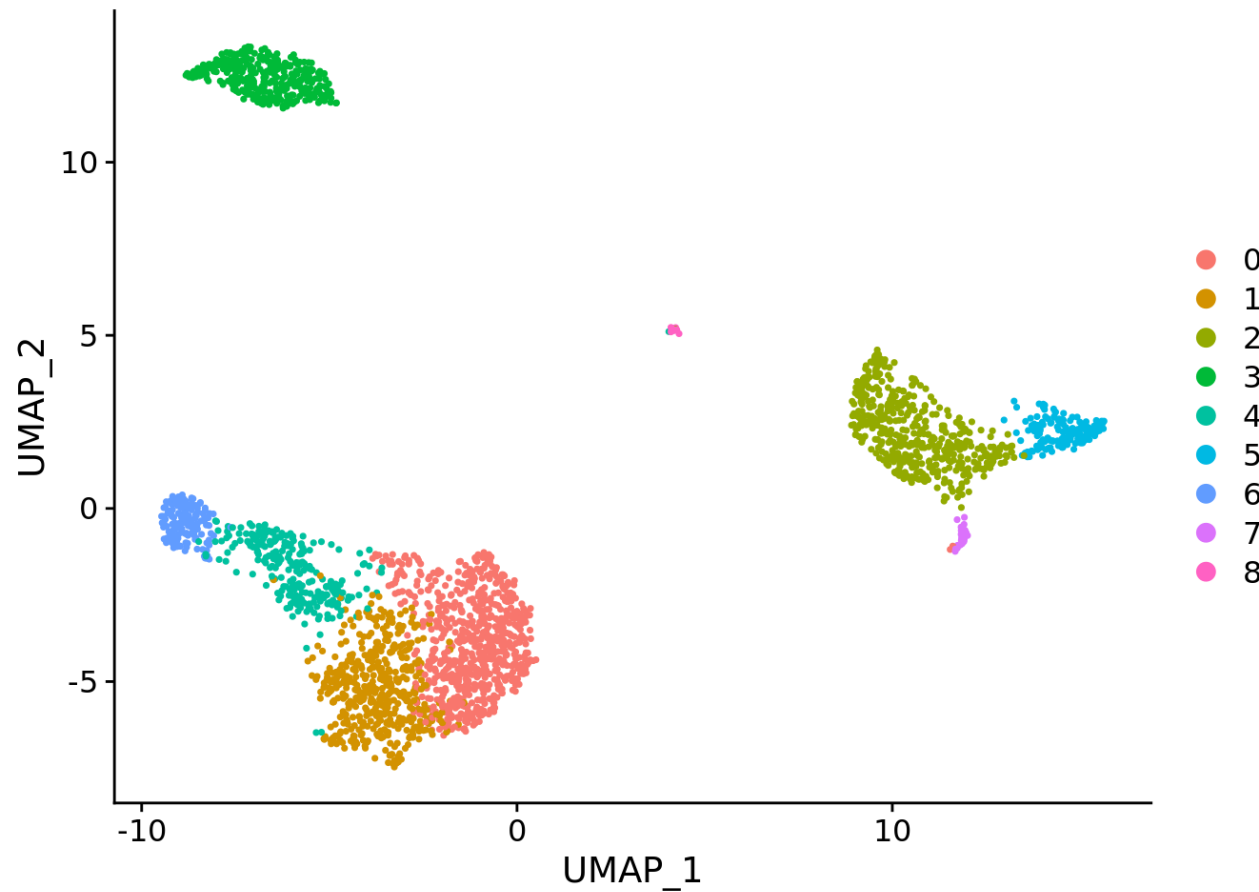
- Find sub-groups in single-cell data using any clustering methods (e.g., Louvain algorithm)

▪ **Assignment step**

- Marker identification: Find differently expressed genes in each cluster compared to other cells (i.e., points in other clusters)
- Assign cell-type identities to the clusters by matching the marker genes with the members in known cell-type signatures

Cell-type identification

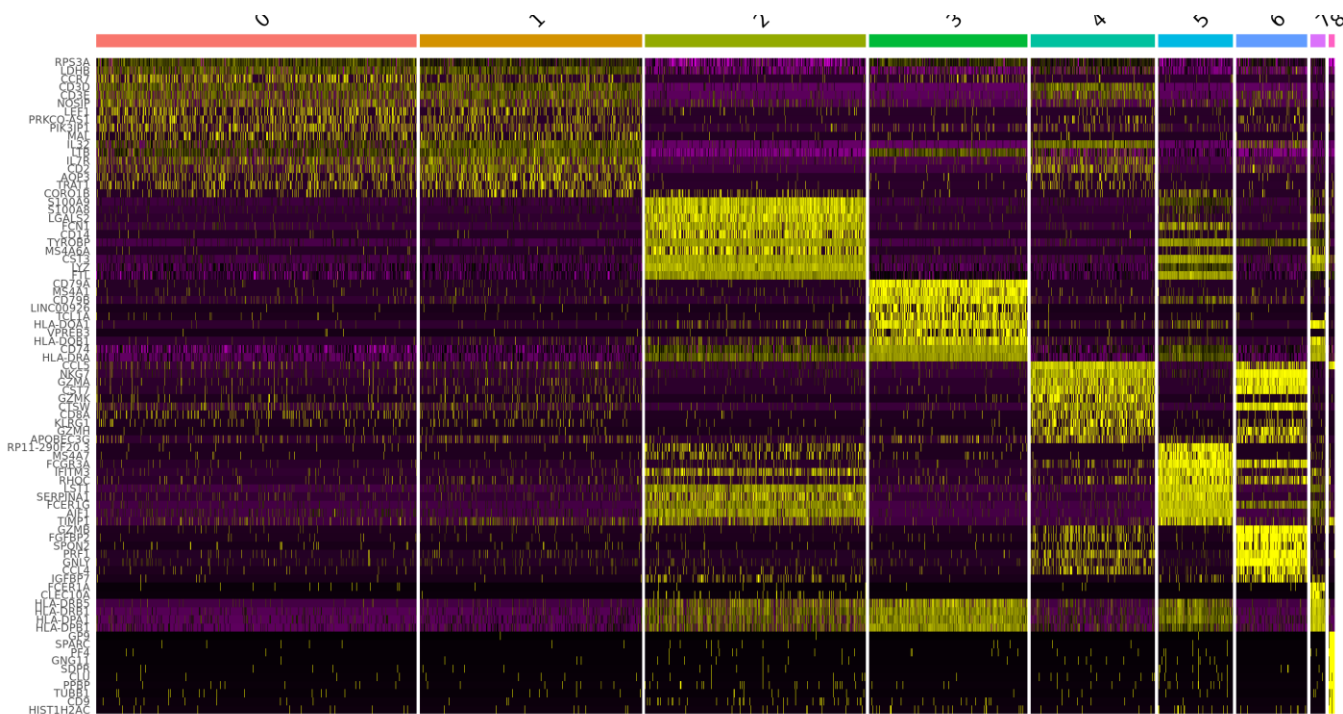
1. Clustering



Cell-type identification

2-1. Marker gene identification: finding differently expressed genes

2-2. Cell-type assignment

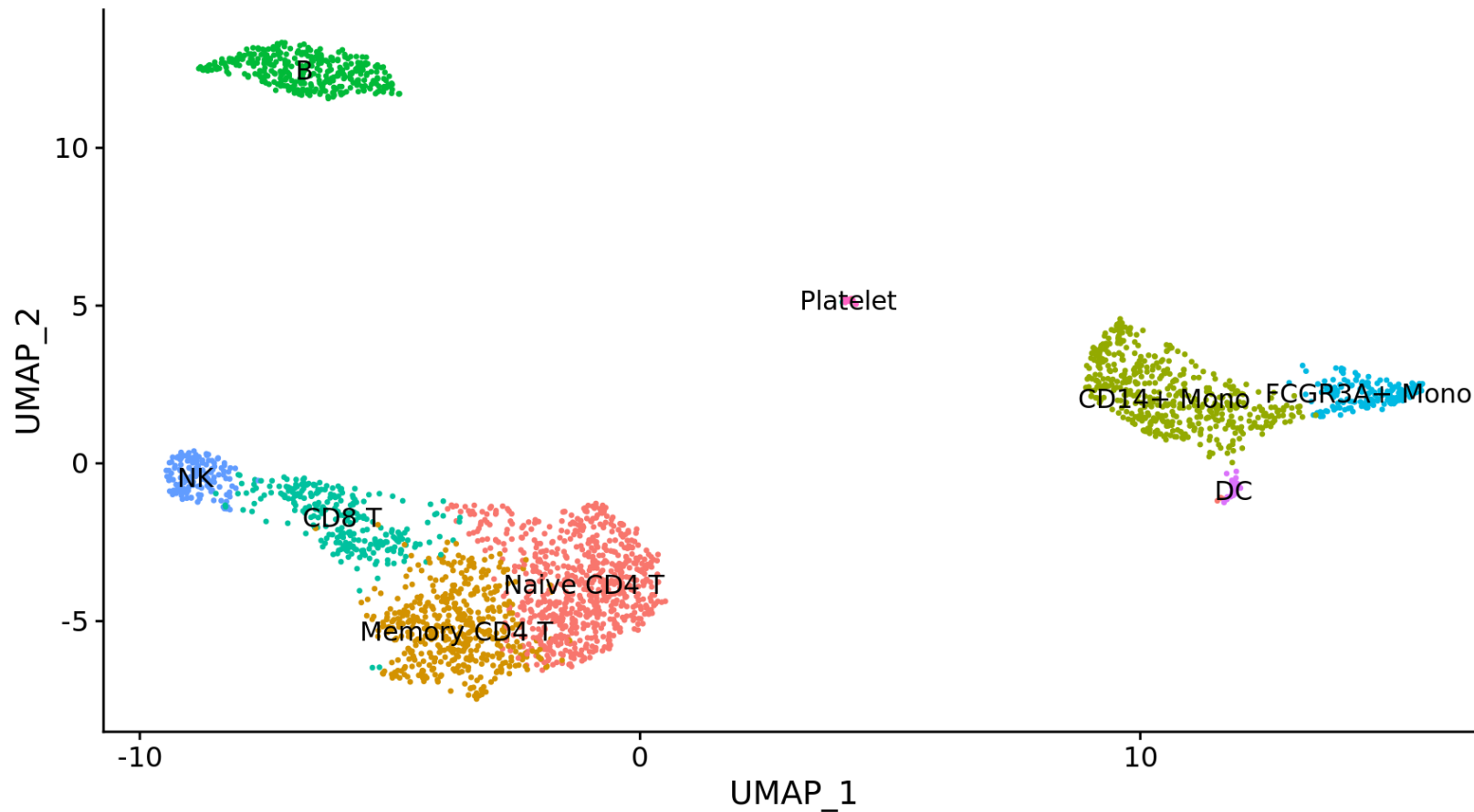


Markers	Cell Type
IL7R, CCR7	Naive CD4+ T
IL7R, S100A4	Memory CD4+
CD14, LYZ	CD14+ Mono
MS4A1	B
CD8A	CD8+ T
FCGR3A, MS4A7	FCGR3A+ Mono
GNLY, NKG7	NK
FCER1A, CST3	DC
PPBP	Platelet

Figure credit: Seurat package

Cell-type identification

- Final results



Outline

- Dimensionality reduction (visualization)
 - t-SNE
 - Uniform Manifold Approximation and Projection (UMAP)
- Clustering approaches
 - Graph based methods: e.g., Louvain algorithm
 - Cell-type identification
- **Multiple Dataset Integration**
 - Canonical correlation analysis & L2-norm normalization
 - Anchoring

Multiple Dataset Integration

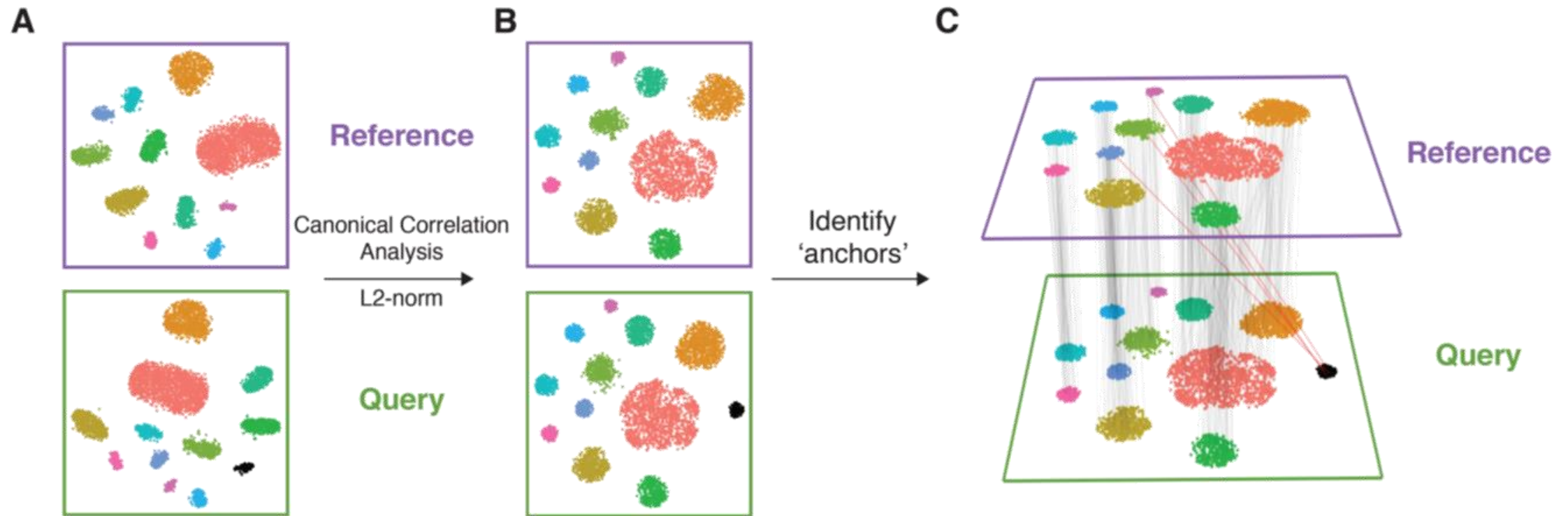


Figure: Tim Stuart et al. "Comprehensive Integration of Single-Cell Data", Cell 2019

Multiple Dataset Integration

- Canonical correlation analysis (CCA)

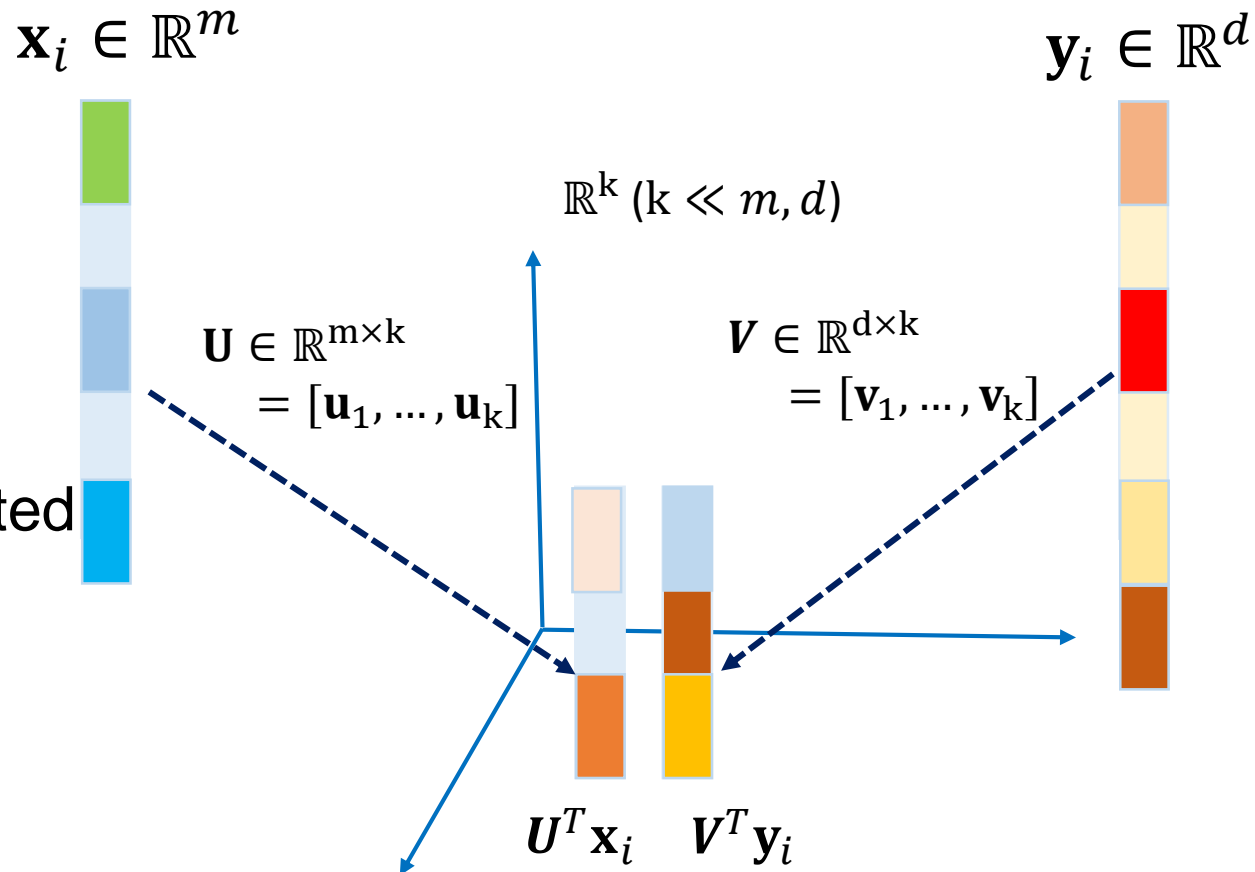
- finds a common space of two sets of data $\mathbf{x}_i \in \mathbb{R}^m$

$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$$

$$\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N]^T$$

- finds a linear transformation of variables that makes the two variables maximally correlated

$$\begin{aligned} \hat{\mathbf{u}}, \hat{\mathbf{v}} &= \operatorname{argmax}_{\mathbf{u}, \mathbf{v}} \mathbf{u}^T \mathbf{X}^T \mathbf{Y} \mathbf{v} \\ \text{s.t. } \mathbf{u}^T \mathbf{X}^T \mathbf{X} \mathbf{u} &\leq 1, \mathbf{v}^T \mathbf{Y}^T \mathbf{Y} \mathbf{v} \leq 1 \end{aligned}$$



Multiple Dataset Integration

- Canonical correlation analysis (CCA)

- finds a common space of two sets of data $\mathbf{x}_i \in \mathbb{R}^m$

$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$$

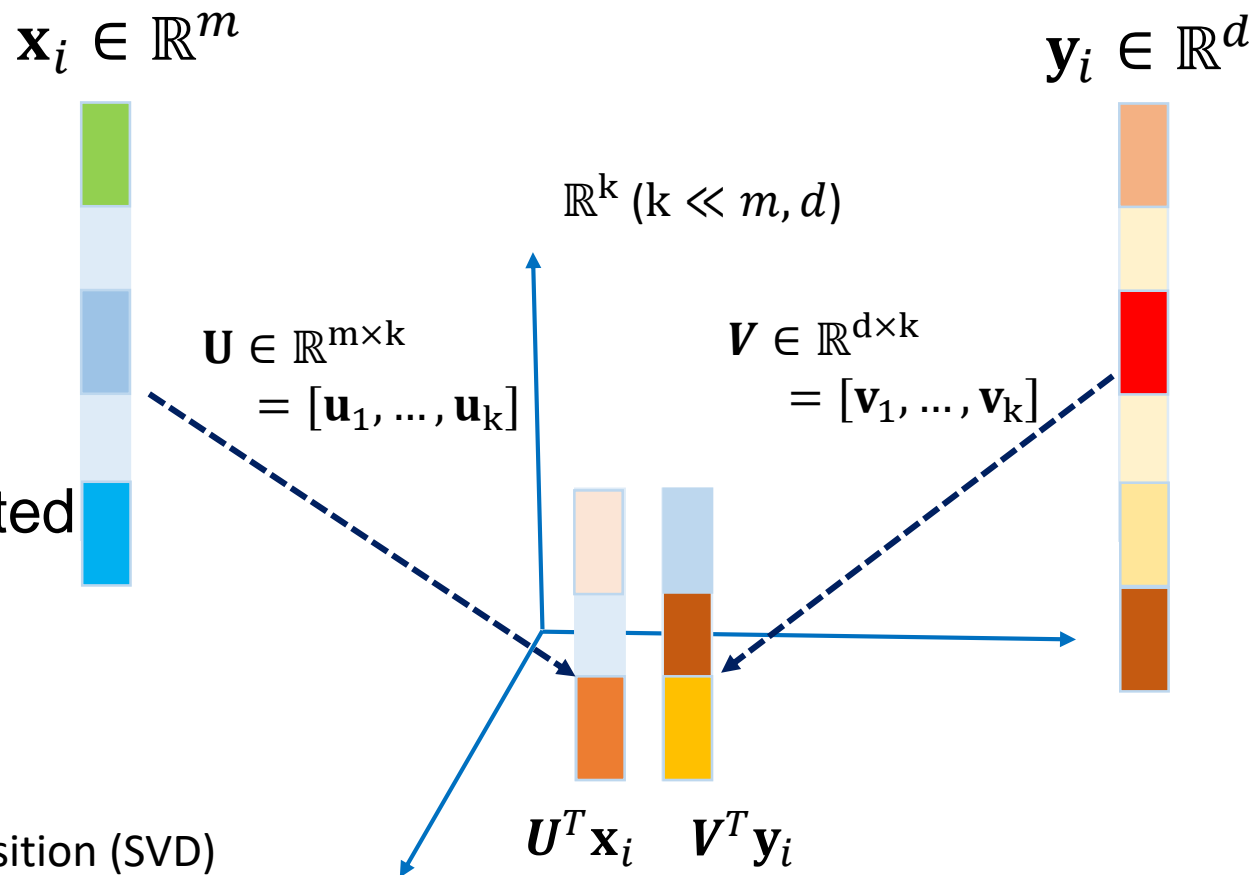
$$\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N]^T$$

- finds a linear transformation of variables that makes the two variables maximally correlated

In Seurat, the maximization problem is simplified as

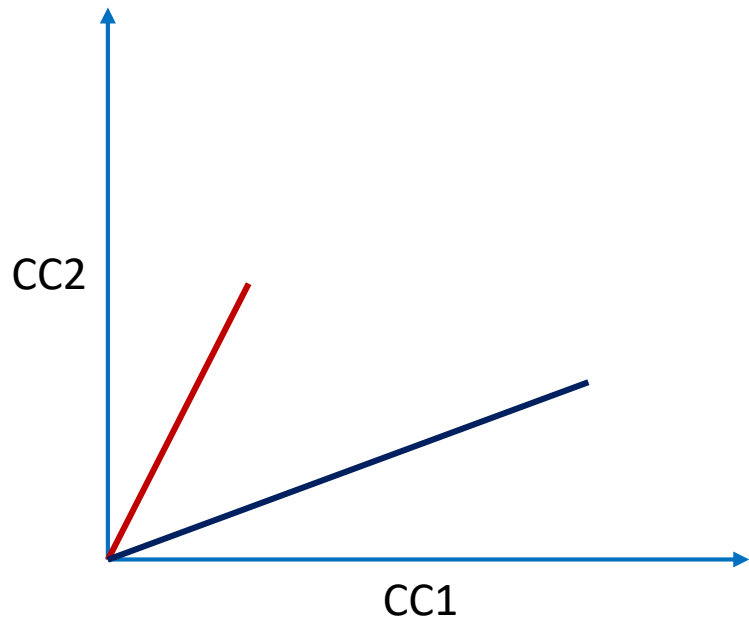
$$\begin{aligned} \hat{\mathbf{u}}, \hat{\mathbf{v}} = \operatorname{argmax}_{\mathbf{u}, \mathbf{v}} & \mathbf{u}^T \mathbf{X}^T \mathbf{Y} \mathbf{v} \\ \text{s.t. } & \|\mathbf{u}\|_2^2 \leq 1, \|\mathbf{v}\|_2^2 \leq 1 \end{aligned}$$

The problem can be solved using singular value decomposition (SVD)

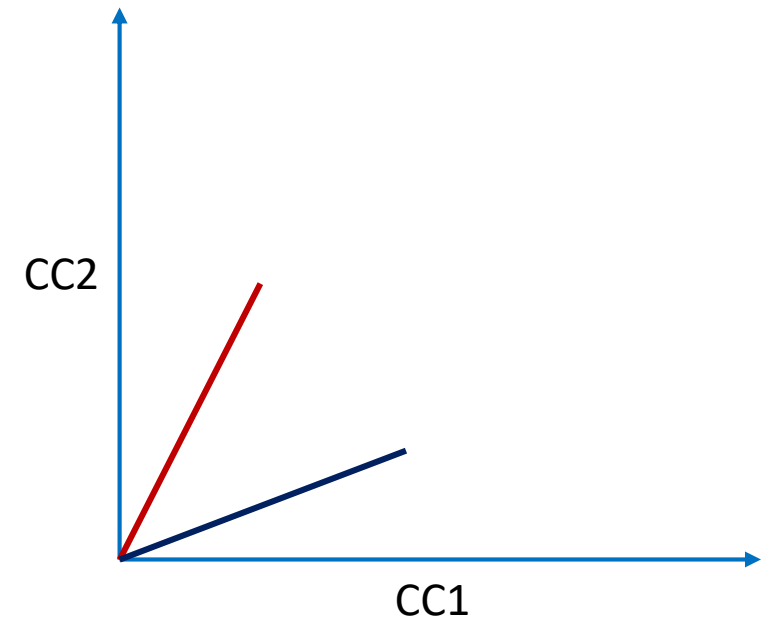


Multiple Dataset Integration

- L_2 -norm normalization

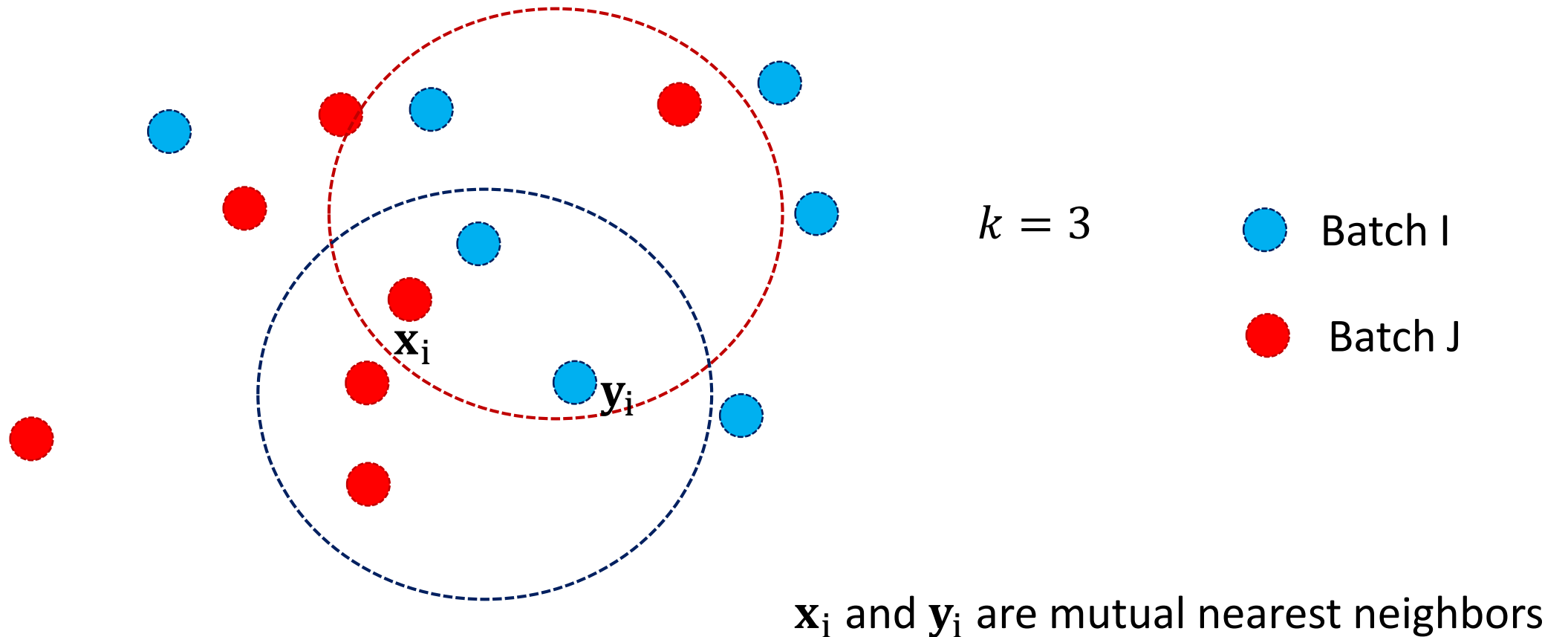


$$\frac{\mathbf{x}}{\|\mathbf{x}\|_2}$$



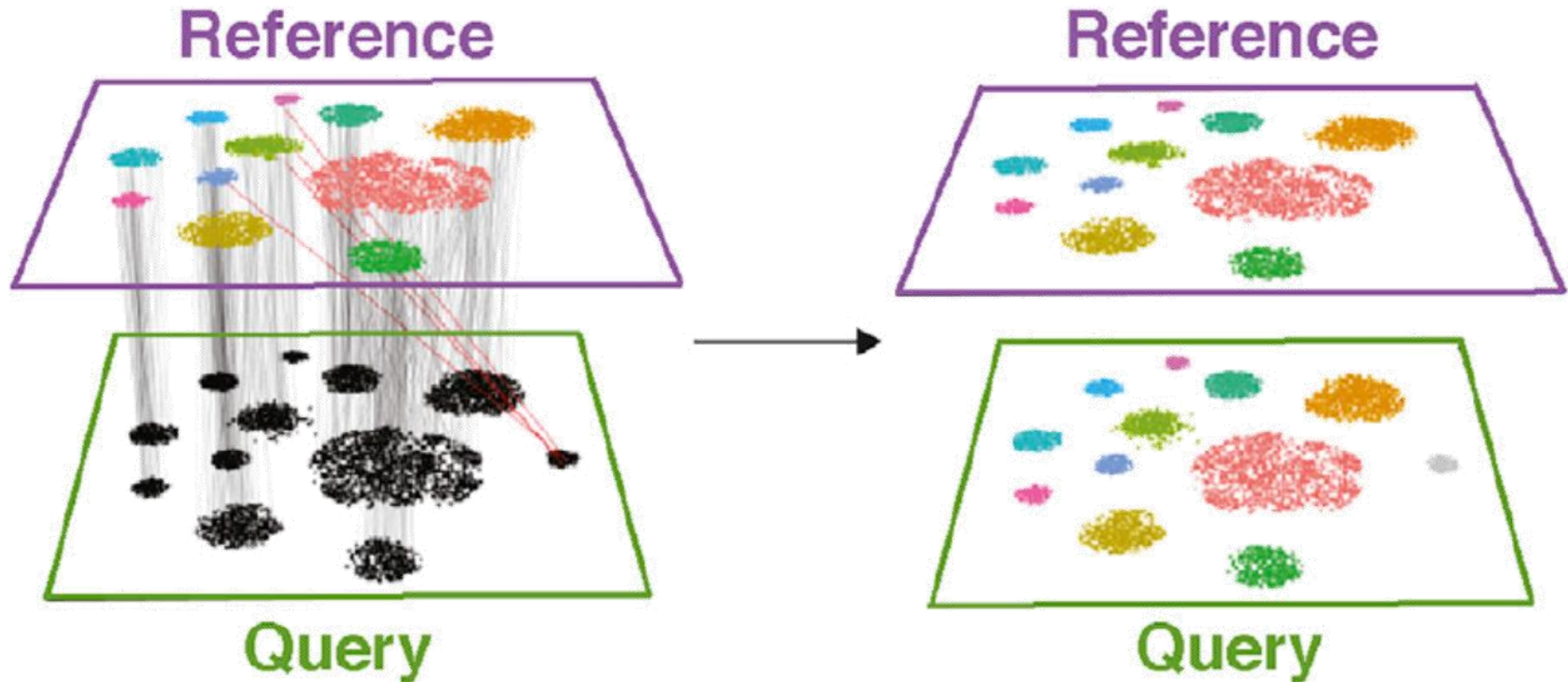
Multiple Dataset Integration

- Anchors (Mutual nearest neighbors)



Multiple Dataset Integration

- Label transfer



References

- Seurat package: <https://satijalab.org/seurat/>
- Stuart, Tim et al. “Comprehensive Integration of Single-Cell Data” Cell, Volume 177, Issue 7, 1888 - 1902.e21, 2019
- How Exactly UMAP Works:
<https://towardsdatascience.com/how-exactly-umap-works-13e3040e1668>
- V. Blondel et al., “Fast unfolding of communities in large networks”, Journal of Statistical Mechanics: Theory and Experiment 2008