# Sparse multiple Canonical Correlation Network discovery (SmCCNet)

Date: 6 Dec 2022

Presenter: Zuqi Li

# About me

**Affiliation:** PhD student at KU Leuven
**Department**: Human Genetics

**Education**: Master in Bioinformatics at the University of Copenhagen
Bachelor in Biology at Beijing Normal University

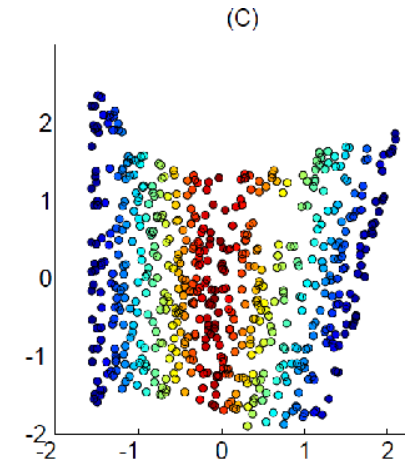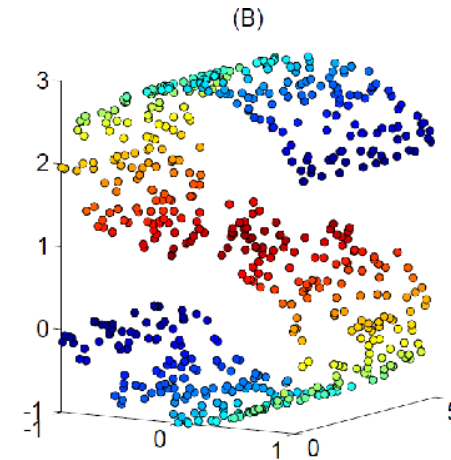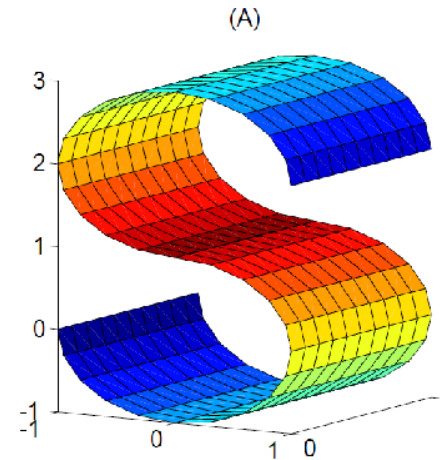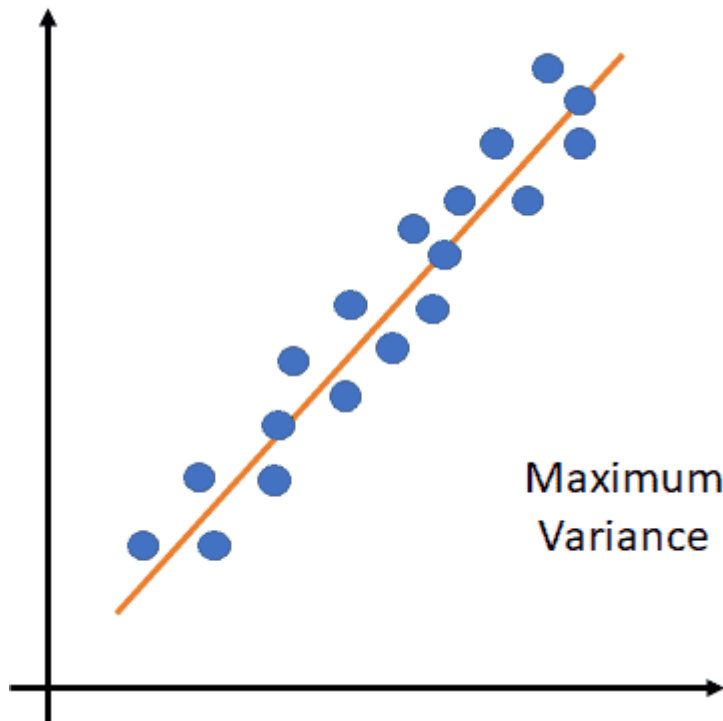**Background**: Machine Learning, computational biology, bioinformatics

**Supervisor**: Prof. Dr. Dr. Kristel Van Steen

Zuqi Li

MARIE CURIE ACTIONS

TRANSYS
PERSONALIZED MEDICINE

# Agenda

- ➤ **PCA**

- ➤ **CCA**

- ➤ **Sparse CCA**

- ➤ **Sparse multiple CCA**

- ➤ **SmCCNet**

- ➤ **Applications**

# Dimensionality reduction

- With the development of technologies and reduction in cost, high-dimensional datasets are being collected.

- Features may be correlated or noisy.



1) https://machinelearninggeek.com/dimensionality-reduction-using-pca/ 2) Engel, Daniel, Lars Hüttenberger, and Bernd Hamann. "A survey of dimension reduction methods for high-dimensional data analysis and visualization." Visualization of Large and Unstructured Data Sets: Applications in Geospatial Planning, Modeling and Engineering-Proceedings of IRTG 1131 Workshop 2011. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2012.
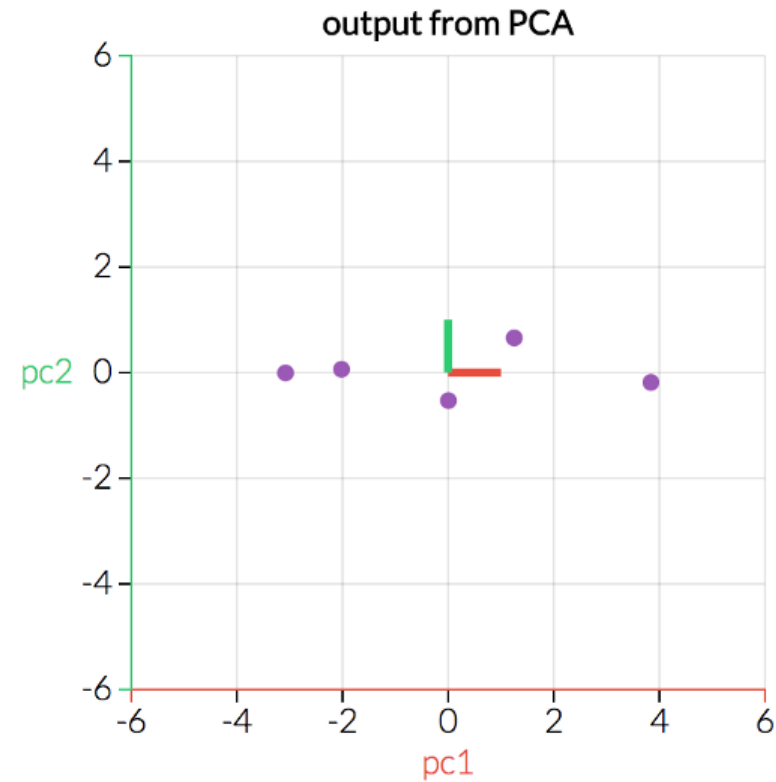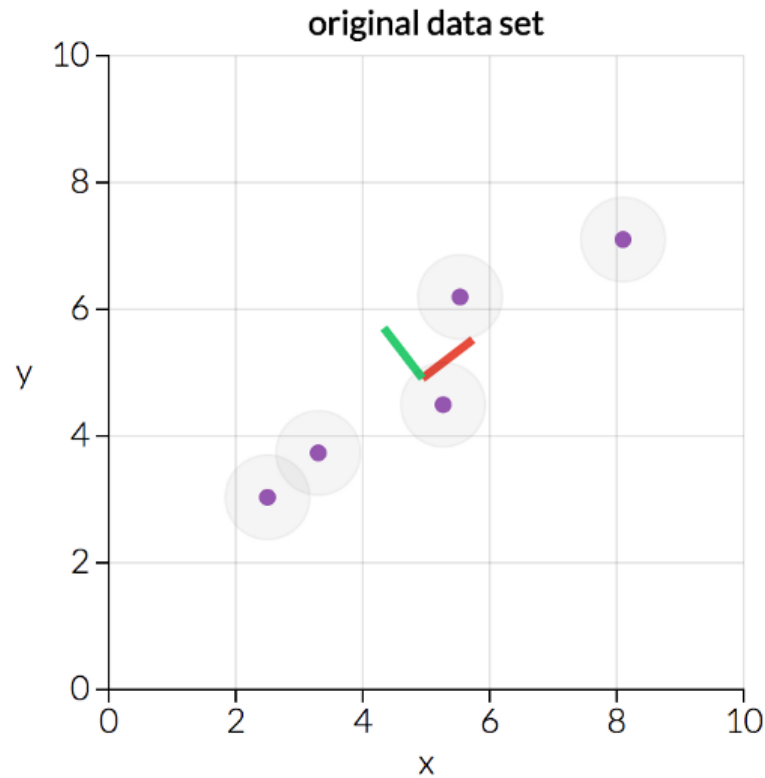
# Principal Component Analysis

- Reduce dimensionality of large datasets and increase interpretability with minimal information loss.

- By creating new uncorrelated variables that successively maximize variance, i.e. principal components.

$$X \in \mathbb{R}^{n \times p} \longrightarrow X^T X \longrightarrow W \in \mathbb{R}^{p \times p} \longrightarrow T = XW$$

$$X = U\Sigma W^T \longrightarrow T = XW = U\Sigma W^T W = U\Sigma$$

# Principal Component Analysis



An example from setosa.io

# Canonical Correlation Analysis

- CCA measures the relatedness of 2 sets of features by maximizing the correlation between some linear transformations of the 2 sets.

- Given 2 standardized data matrices $X_1 \in R^{n \times p_1}$ and $X_2 \in R^{n \times p_2}$, their canonical correlation is $Corr(X_1 w_1, X_2 w_2) = w_1^T X_1^T X_2 w_2$.
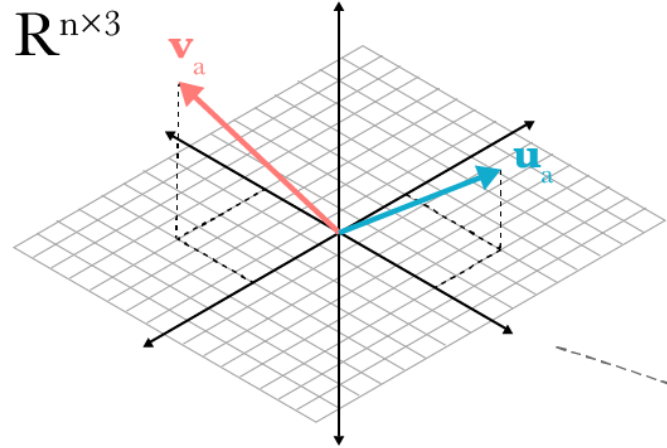
$$(w_1, w_2) = \arg\max_{\widetilde{w}_1, \widetilde{w}_2} (\widetilde{w}_1^T X_1^T X_2 \widetilde{w}_2) \ \ s.t. \|\widetilde{w}_1\|^2 = \|\widetilde{w}_2\|^2 = 1$$

$$u = X_1 w_1$$
$$v = X_2 w_2$$

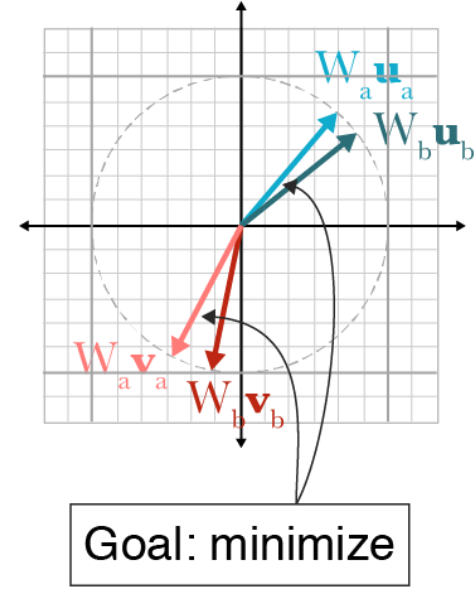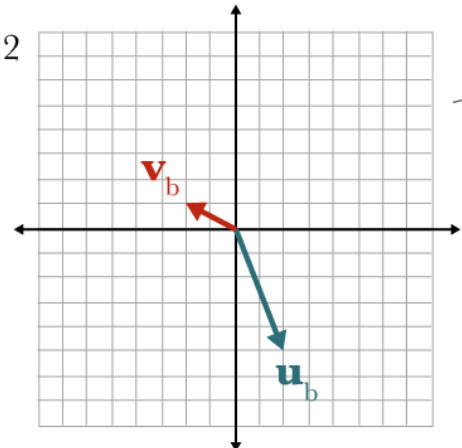# Canonical Correlation Analysis



$X_a \in R^{n \times 3}$

$\mathbf{v}_a$

$\mathbf{u}_a$

$X_a W_a, X_b W_b \in R^{n \times 2}$

$W_a \in R^{3 \times 2}$

$X_b \in R^{n \times 2}$

$\mathbf{v}_b$

$\mathbf{u}_b$

$W_b \in R^{2 \times 2}$

$W_a \mathbf{u}_a$

$W_b \mathbf{u}_b$

$W_a \mathbf{v}_a$

$W_b \mathbf{v}_b$

Goal: minimize

MARIE CURIE ACTIONS

TRANSYS
PERSONALIZED MEDICINE

# Sparse CCA

- In practice, not all features contribute to the overall canonical correlation. Or we may want to reduce the feature set.

- Sparsity is imposed to the canonical weights by adding convex penalty functions $P_1(\cdot), P_2(\cdot)$.

$$(w_1, w_2) = \arg\max_{\widetilde{w}_1, \widetilde{w}_2} (\widetilde{w}_1^T X_1^T X_2 \widetilde{w}_2)$$

$$s.t. \|\widetilde{w}_1\|^2 = \|\widetilde{w}_2\|^2 = 1,$$

$$P_1(\widetilde{w}_1) \leq c_1, P_2(\widetilde{w}_2) \leq c_2$$

# Sparse multiple CCA

- Apart from $X_1$, $X_2$, we may want to take an extra dataset $X_3$ into account of the canonical correlation.

- Use coefficients to prioritize the pairwise correlations.

$$(w_1, w_2, w_3) = \underset{\widetilde{w}_1, \widetilde{w}_2, \widetilde{w}_3}{\arg \max} (a\widetilde{w}_1^T X_1^T X_2 \widetilde{w}_2 + b\widetilde{w}_1^T X_1^T X_3 \widetilde{w}_3 + c\widetilde{w}_2^T X_2^T X_3 \widetilde{w}_3)$$

$$s.t. \|\widetilde{w}_s\|^2 = 1, P_s(\widetilde{w}_s) \leq c_s, s = 1,2,3$$

# Sparse multiple CCA (2)

- Apart from $X_1, X_2$, there's also a phenotype of interest $Y$ that has been measured for the same $n$ subjects.

- The 3rd dataset is a phenotype, i.e. column vector.

$$(w_1, w_2) = \arg\max_{\widetilde{w}_1, \widetilde{w}_2} \left( a\widetilde{w}_1^T X_1^T X_2 \widetilde{w}_2 + b\widetilde{w}_1^T X_1^T Y + c\widetilde{w}_2^T X_2^T Y \right)$$

$$s.t. \|\widetilde{w}_s\|^2 = 1, P_s(\widetilde{w}_s) \leq c_s, s = 1,2$$

PCA

$\downarrow$

CCA

$\downarrow$

Sparse CCA

$\downarrow$

Sparse multiple CCA

$\downarrow$

SmCCNet

# Unsupervised discovery of phenotype specific multi-omics networks

*W Jenny Shi, Yonghua Zhuang, et al.*

# Background and Motivation

- Different quantitative omics measurements on the same subjects.
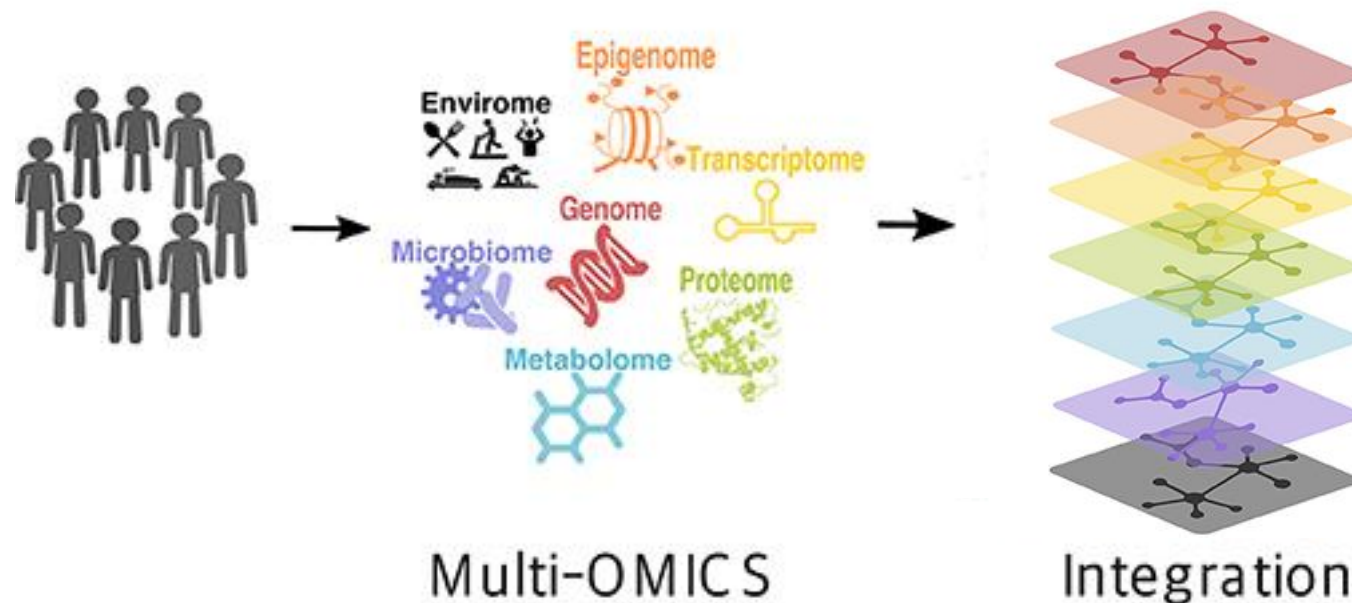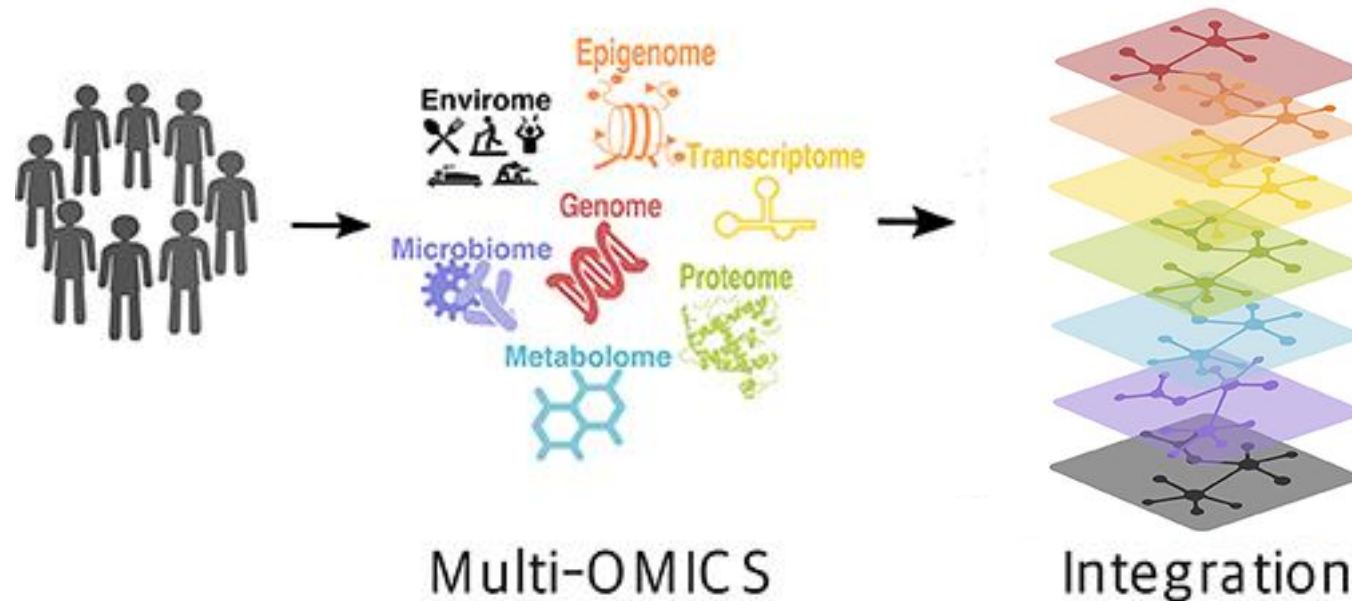- Combine multiple omics data types to study complex traits.
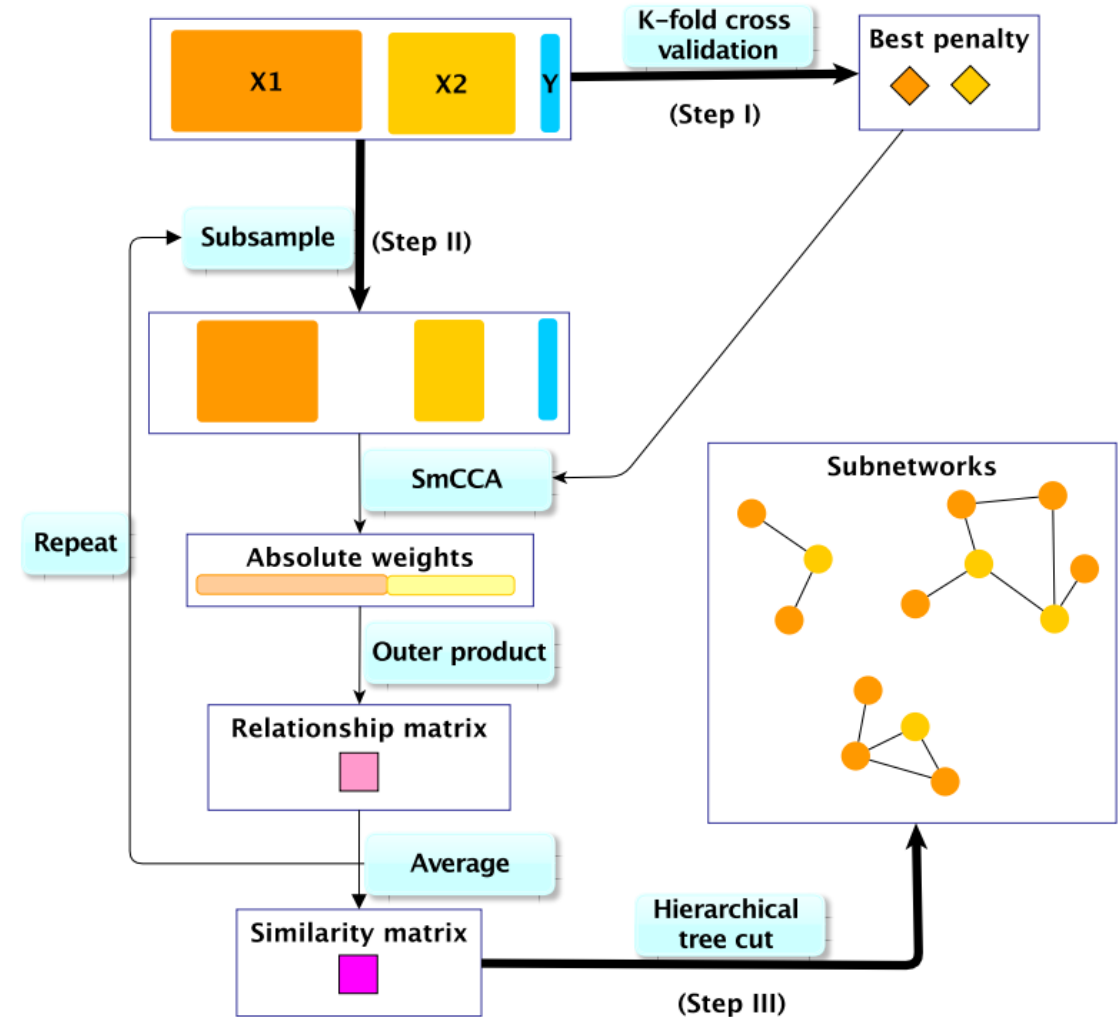
# Background and Motivation



- Supervised methods identify features that are most predictive of the phenotype.

- Inform the integrated feature network with phenotpye.

# Sparse multiple Canonical Correlation Network discovery (SmCCNet)

Simultaneously integrates multiple omics profiles and phenotype information to build interpretable network that models the underlying mechanisms.
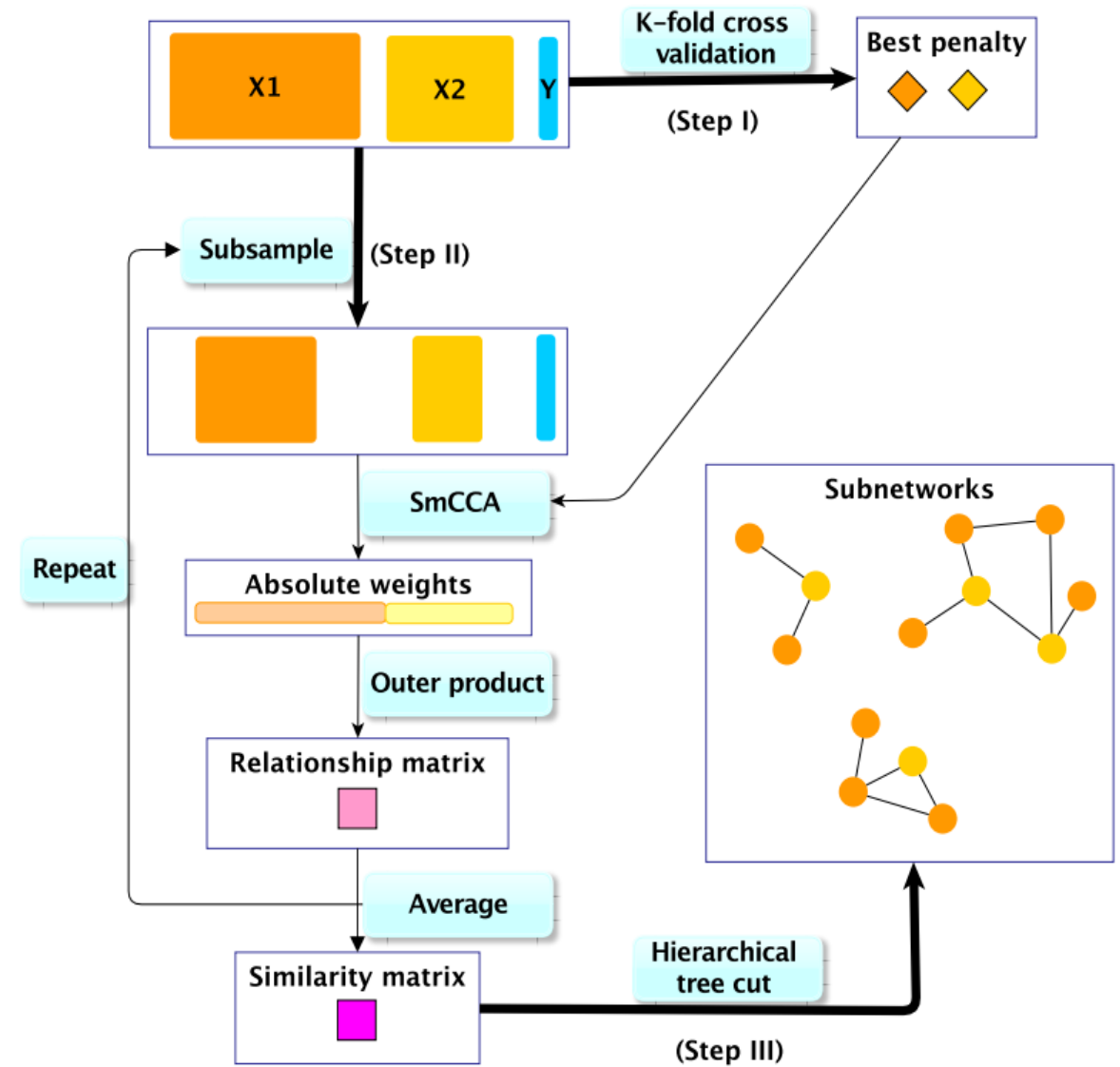
- Sparsity
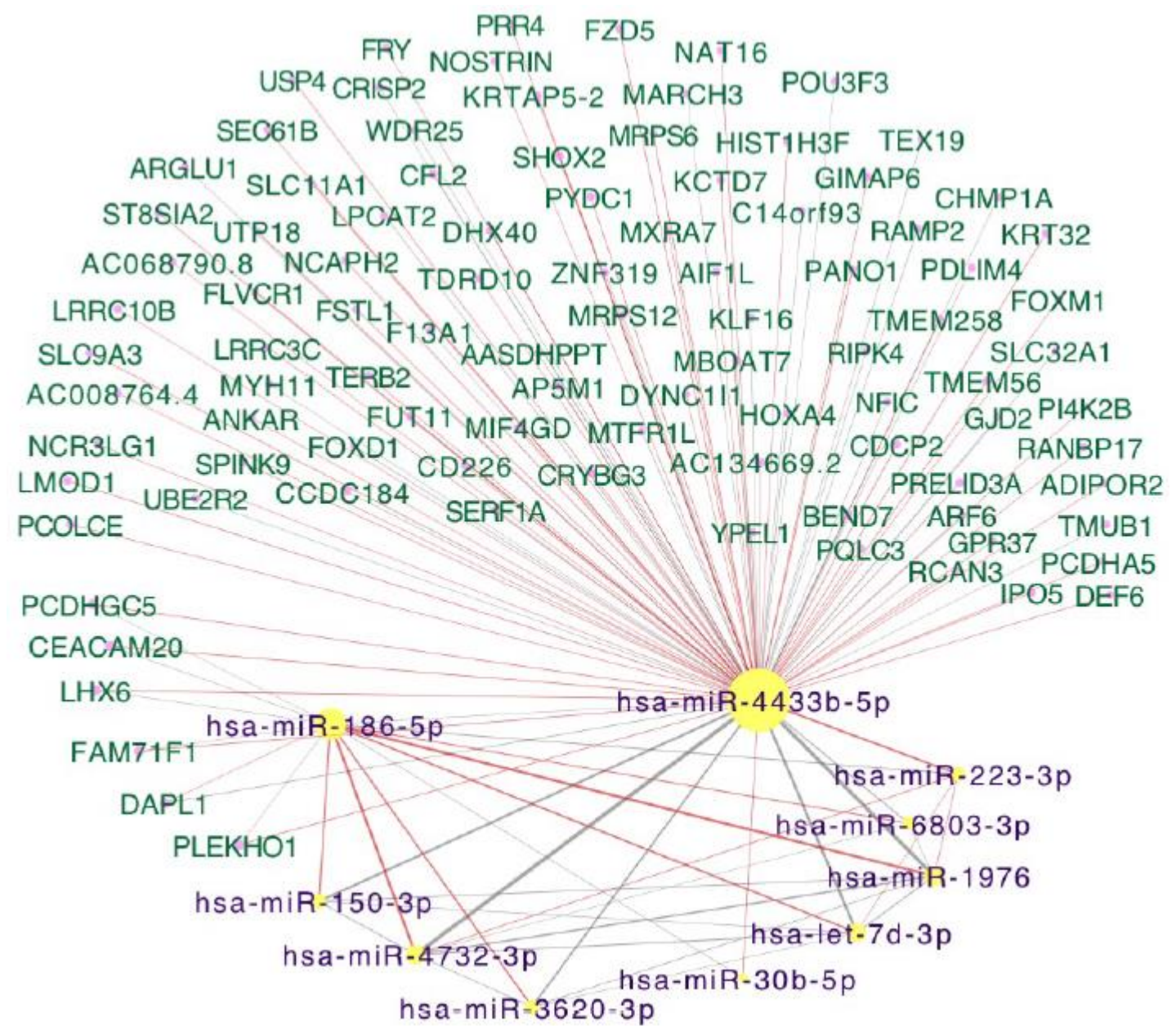- Additional trait
- Module detection

# SmCCNet workflow

1. Identify the best sparsity penalty parameters.
2. Generate robust canonical weights
   A. Feature subsampling
   B. Relationship matrices
   C. Similarity matrix
3. Hierarchical tree cut
   A. Modules
   B. Edge threshold

# Evaluation: COPD data

- 27 subjects (13 controls and 14 cases)
  - 414 miRNAs and 5001 mRNAs
  - Forced expiratory volume during the first second
- 12 connected miRNA-mRNA modules
  - 14,694 negative connections
  - 147 miRNA-mRNA targets have been validated
  - 988 additional targets have been predicted using MultiMir
- SmCCNet identified a higher percentage of predicted and validated miRNA-mRNA pairs than SsCCA
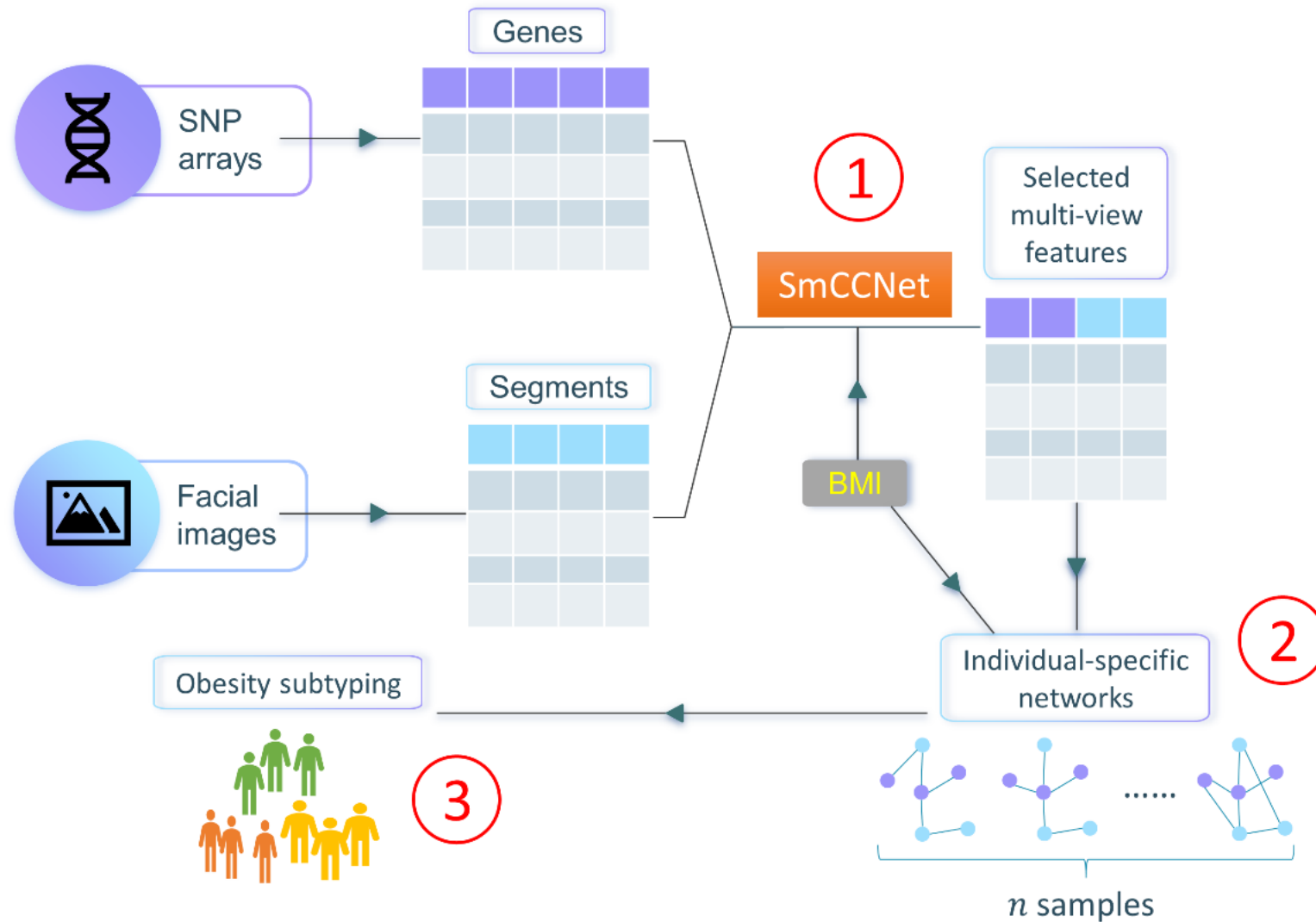
- After edge thresholding, only module 2 remains (10 miRNA and 97 genes)
- miR-4433b-5p is a hub, connected to all 97 genes and 10 miRNAs
  - A biomarker for multidrug-resistance tuberculosis.
- miR-186-5p has been found to be up-regulated in COPD patients.

# Applications of SmCCNet

- SmCCNet can effectively construct phenotype-specific multi-omics network and detect informative modules.

- Features contained in such informed modules can be used for subsequent analyses.

- "A novel network-guided multi-view clustering workflow: dissecting genetic and facial heterogeneity" (netMUG)
    - Data: genomics and facial images
    - Trait: BMI
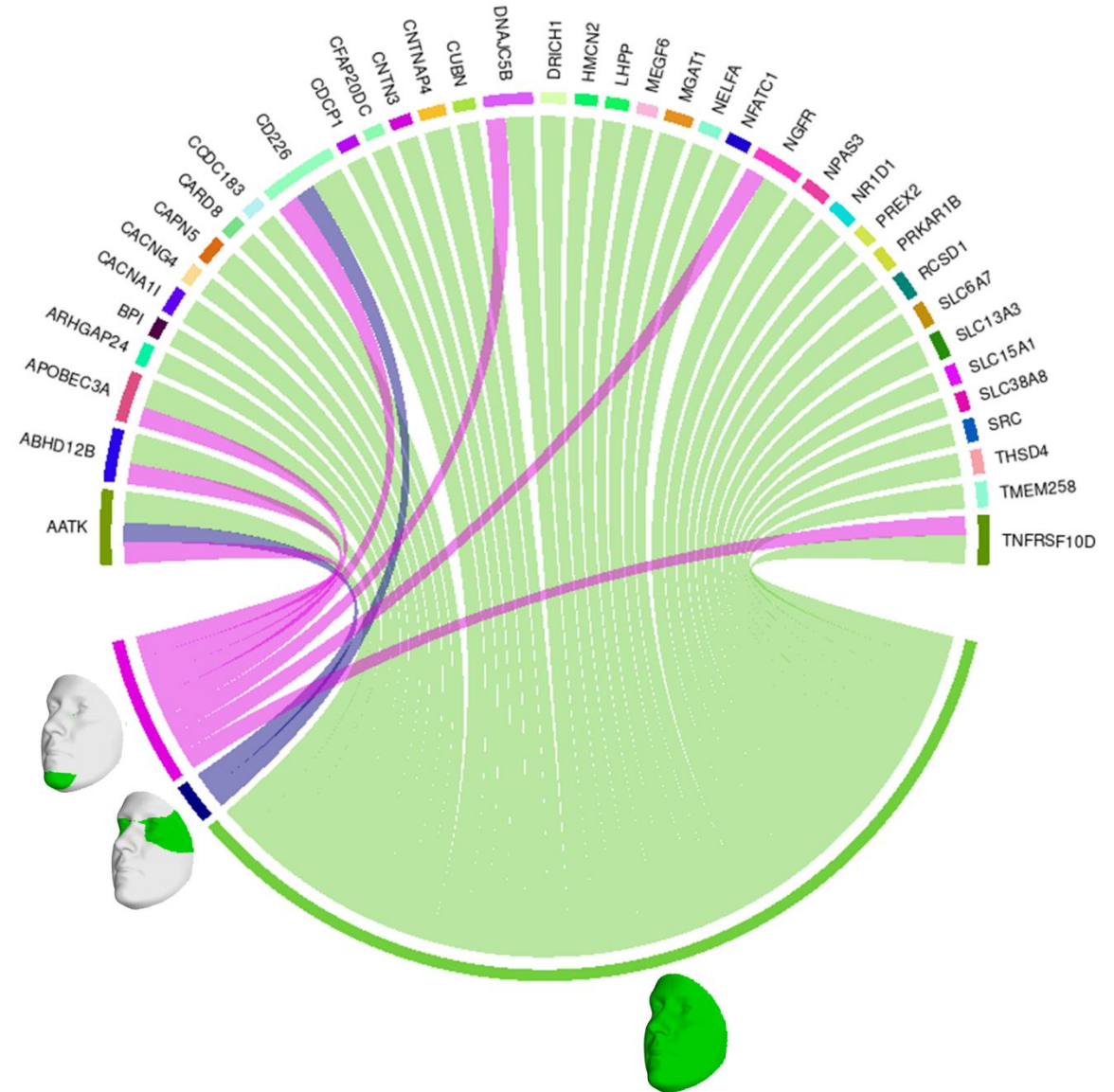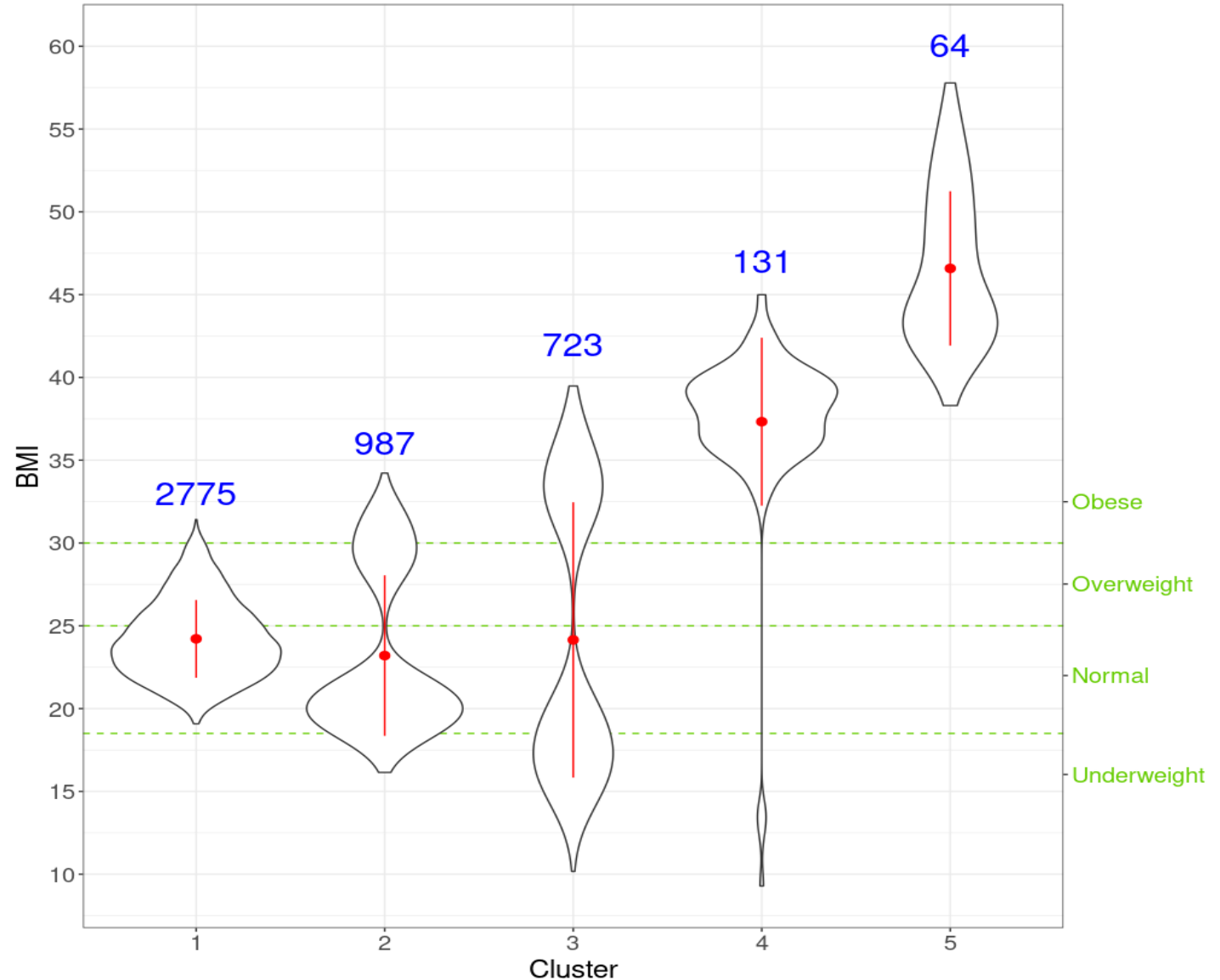    - Goal: obesity sutyping

# netMUG workflow

# Features selected by SmCCNet

- SmCCNet selected 278 genes and 26 facial segments.
  - 150 genes are also found by GWAS
  - 39 genes are also highlighted in DisGeNET
- Top 1% connections in the SmCCNet network:
  - AATK is known highly associated with BMI
  - APOBEC3A, DNAJC5B and NGFR all affect body height

# Obesity subtyping

- netMUG found 5 clusters
- The clusters are significantly associated with BMI
- Our subtyping is complementary to the classic BMI categories.

# Obesity subtyping

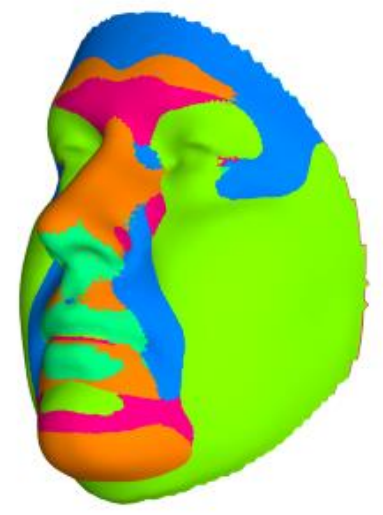

Cluster 1   Cluster 2   Cluster 3   Cluster 4   Cluster 5   All clusters

# Take-home messages

- PCA is a dimensionality reduction method for a single dataset.
- CCA works with two datasets and maximizes their canonical correlation.
- Sparse CCA prioritizes features most contributing to the correlation.
- Sparse multiple CCA can incorporate more than two datasets.
- SmCCNet constructs phenotype-specific multi-omic network and can be used as feature selector.

# Thank you!
# Q&A