

Inference after human genetic clustering

Javier González-Delgado

University of Rennes, ENSAI, CREST

MASCOT-NUM 2026

April 2, 2026



Genetic diversity and population structure

SNP data describing genotypes in a population

		INDIVIDUAL					
{	GENOTYPES	AA	AT-AC	CC	GG	AC	1
		AG	AT-AA	CC	GT	CC	2
		AG	TT-AA	AA	GT	CC	3
		AG	TT-AA	AC	GG	AA	4

SNP data describing genotypes in a population

		INDIVIDUAL					
GENOTYPES	}	AA	AT-AC	CC	GG	AC	1
		AG	AT-AA	CC	GT	CC	2
		AG	TT-AA	AA	GT	CC	3
		AG	TT-AA	AC	GG	AA	4

Allele configurations → SNP code

AA → 0, Aa → 1, aa → 2,

with A (resp. a) being the reference (resp. alternate) allele.

SNP data describing genotypes in a population

		INDIVIDUAL										
{	AA	—	AT	—	AC	—	CC	—	GG	—	AC	1
	AG	—	AT	—	AA	—	CC	—	GT	—	CC	2
	AG	—	TT	—	AA	—	AA	—	GT	—	CC	3
	AG	—	TT	—	AA	—	AC	—	GG	—	AA	4

Allele configurations → SNP code

AA → 0, Aa → 1, aa → 2,

with A (resp. a) being the reference (resp. alternate) allele.

		SNPs									
		1	2	3	4	5	6	7	8	9	10
{	1	0	0	2	2	1	0	0	0	2	2
	2	2	1	0	1	2	1	2	1	1	1
	3	1	2	0	1	0	2	2	1	0	0
	4	0	1	2	2	2	0	1	0	2	2
	5	1	2	0	0	0	1	2	1	1	0
	6	1	0	2	1	2	0	0	1	1	2
	7	0	0	1	2	2	1	0	0	2	1
	8	2	2	0	0	1	2	0	2	0	1

Figures : *An introduction to human population genetics, variation, and disease*, Jonathan K Pritchard.

SNP data describing genotypes in a population

		INDIVIDUAL										
{	AA	—	AT	—	AC	—	CC	—	GG	—	AC	1
	AG	—	AT	—	AA	—	CC	—	GT	—	CC	2
	AG	—	TT	—	AA	—	AA	—	GT	—	CC	3
	AG	—	TT	—	AA	—	AC	—	GG	—	AA	4

Allele configurations → SNP code

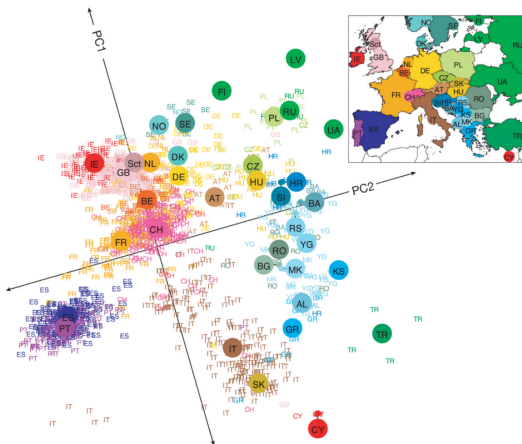
AA → 0, Aa → 1, aa → 2,

with A (resp. a) being the reference (resp. alternate) allele.

		SNPs									
		1	2	3	4	5	6	7	8	9	10
INDIVIDUALS	1	0	0	2	2	1	0	0	0	2	2
	2	2	1	1	0	1	2	1	2	1	1
	3	1	2	0	1	0	2	2	1	0	0
	4	0	1	2	2	2	0	1	0	2	2
	5	1	2	0	0	0	1	2	1	1	0
	6	1	0	2	1	2	0	0	1	1	2
	7	0	0	1	2	2	1	0	0	2	1
	8	2	2	0	0	1	2	0	2	0	1

→ Is it possible to find structure in that genotype matrix?

Genotype data unveils population structure



PCA on genome-wide SNP data for 1387 European individuals (Novembre *et al.*, Nature, 2008).

Population structure and genotype data

Fundamental question in Population Genetics

Can we infer/observe population structure and/or ancestry from genetic data?

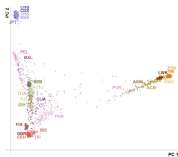
Population structure and genotype data

Fundamental question in Population Genetics

Can we infer/observe population structure and/or ancestry from genetic data ?

INDIVIDUALS	SNPs									
	1	2	3	4	5	6	7	8	9	10
1	0	0	2	2	1	0	0	0	2	2
2	2	1	1	0	1	2	1	2	1	1
3	1	2	0	1	0	2	2	1	0	0
4	0	1	2	2	2	0	1	0	2	2
5	1	2	0	0	0	1	2	1	1	0
6	1	0	2	1	2	0	0	1	1	2
7	0	0	1	2	2	1	0	0	2	1
8	2	2	0	0	1	2	0	2	0	1

PCA



Population
structure

Population structure and genotype data

Fundamental question in Population Genetics

Can we infer/observe population structure and/or ancestry from genetic data?

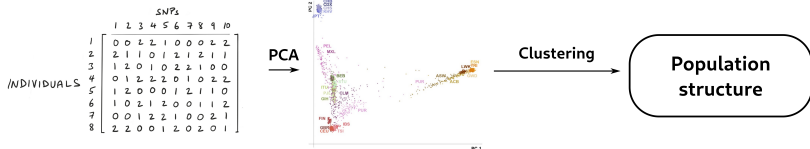
INDIVIDUALS	SNPs									
	1	2	3	4	5	6	7	8	9	10
1	0	0	2	2	1	0	0	0	2	2
2	2	1	1	0	1	2	1	2	1	1
3	1	2	0	1	0	2	2	1	0	0
4	0	1	2	2	2	0	1	0	2	2
5	1	2	0	0	0	1	2	1	1	0
6	1	0	2	1	2	0	0	1	1	2
7	0	0	1	2	2	1	0	0	2	1
8	2	2	0	0	1	2	0	2	0	1



Population structure and genotype data

Fundamental question in Population Genetics

Can we infer/observe population structure and/or ancestry from genetic data?



- Simple clustering algorithms like hierarchical clustering reveal broad (continental) structure¹.

1. Mountain and Cavalli-Sforza (1997).

Population structure and genotype data

Fundamental question in Population Genetics

Can we infer/observe population structure and/or ancestry from genetic data ?



- Simple clustering algorithms like hierarchical clustering reveal broad (continental) structure¹.
- More complex algorithms can reveal population structure at finer levels² (e.g. the UMAP+HDBSCAN($\hat{\epsilon}$) pipeline³).

1. Mountain and Cavalli-Sforza (1997), 2. Rosenberg *et al.* (2002), Hoffelder *et al.* (2017), 3. Diaz-Papkovich *et al.* (2025).

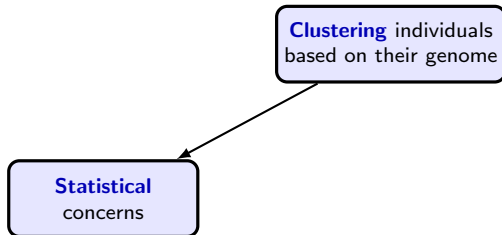
Clustering genetic data

Main challenges

Clustering individuals
based on their genome

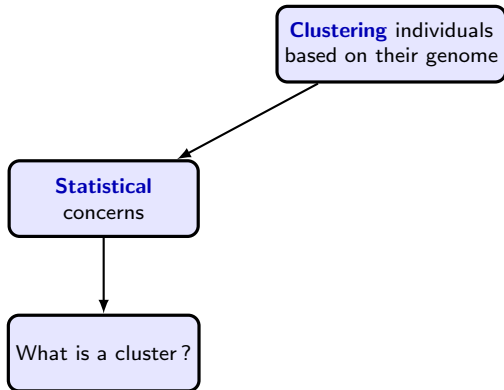
Clustering genetic data

Main challenges



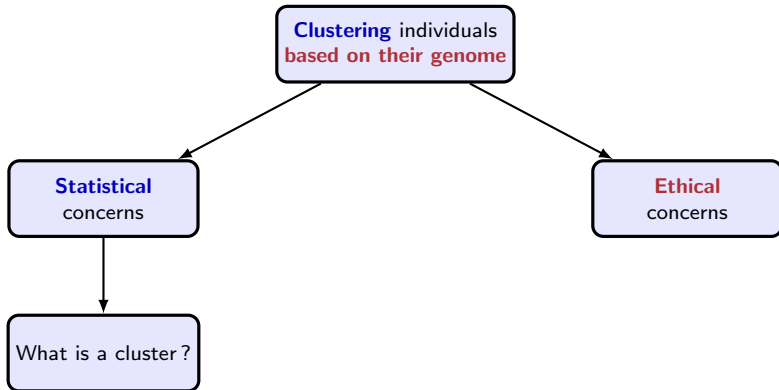
Clustering genetic data

Main challenges



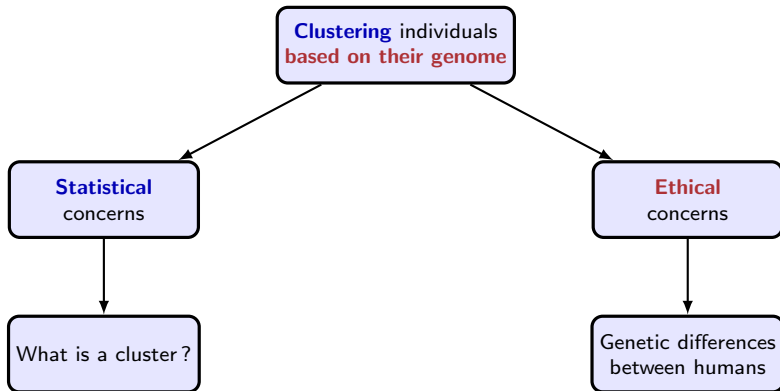
Clustering genetic data

Main challenges



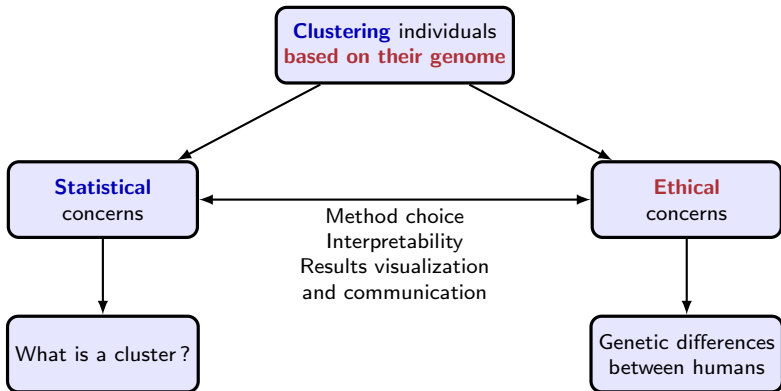
Clustering genetic data

Main challenges



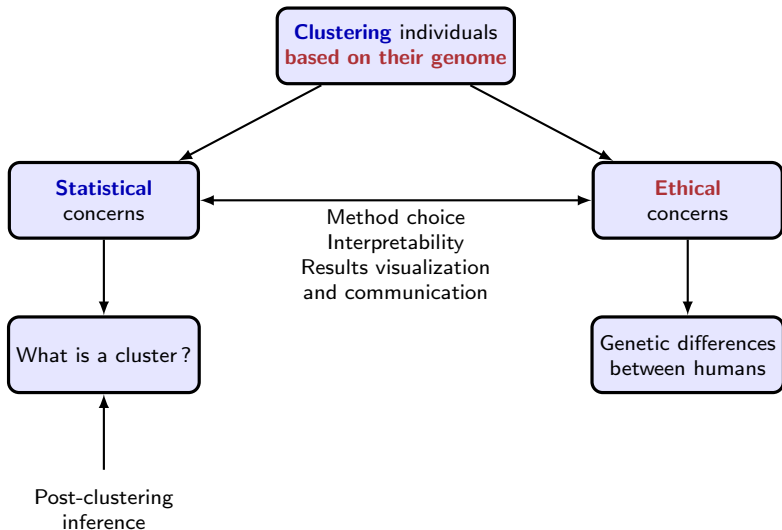
Clustering genetic data

Main challenges



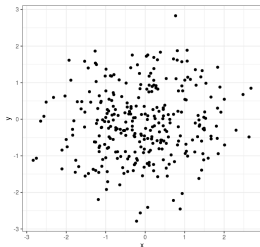
Clustering genetic data

Main challenges



Post-clustering inference

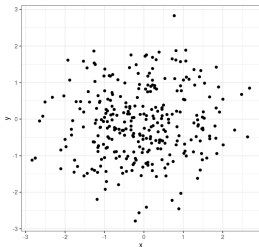
The problem of *double-dipping*



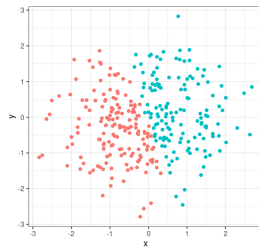
Data

Post-clustering inference

The problem of *double-dipping*



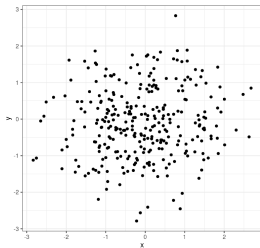
Data



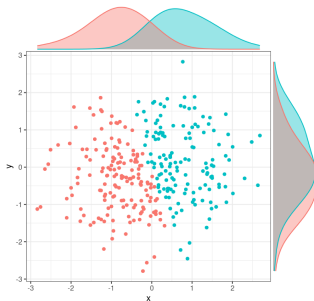
Clustering

Post-clustering inference

The problem of *double-dipping*



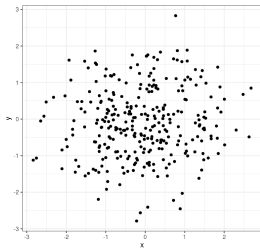
Data



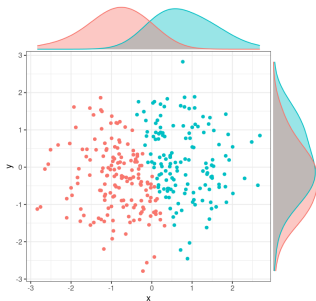
Cluster ● = Cluster ●?

Post-clustering inference

The problem of *double-dipping*



Data



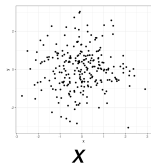
Cluster ● = Cluster ●?

Double dipping

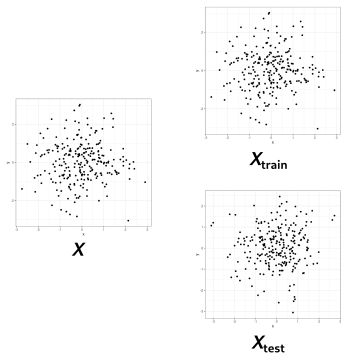
Same data used for both clustering and testing \Rightarrow Lost of type I error control.

Post-clustering inference approaches

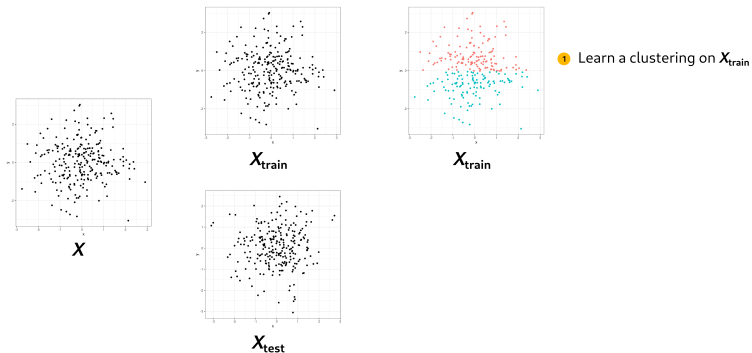
Sample splitting is not a solution



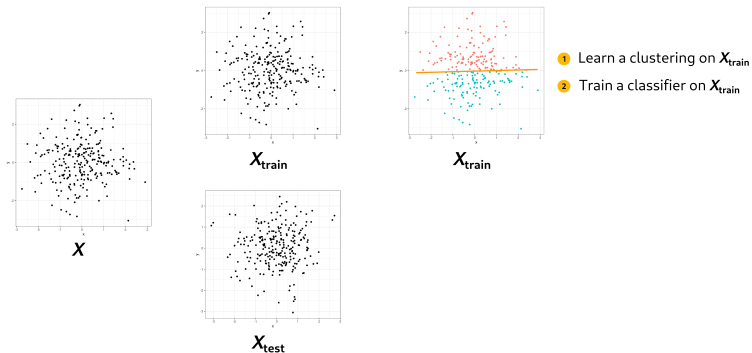
Sample splitting is not a solution



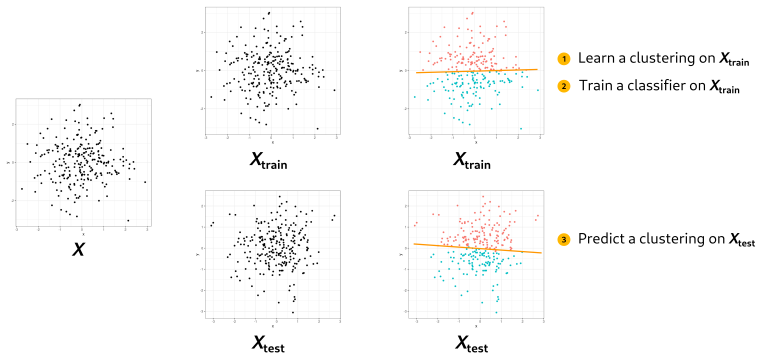
Sample splitting is not a solution



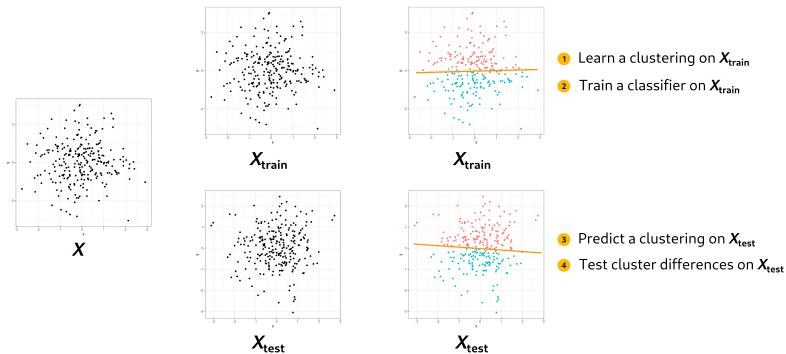
Sample splitting is not a solution



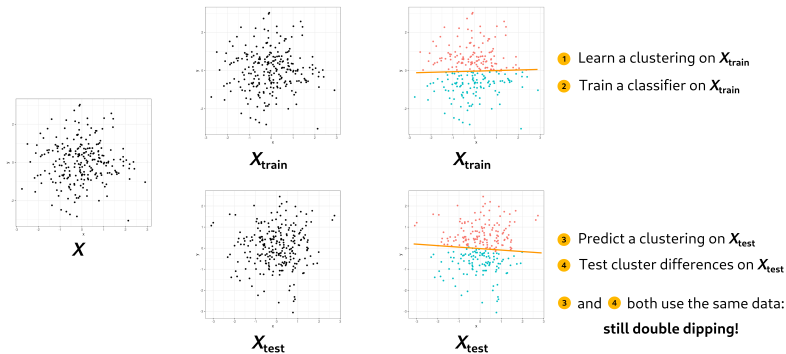
Sample splitting is not a solution



Sample splitting is not a solution



Sample splitting is not a solution



Information partitioning approaches

Data fission / data thinning

Information partitioning approaches

Data fission / data thinning

Idea

Decompose the dataset \mathbf{X} into **two independent datasets** $\mathbf{X}^{(1)}$, $\mathbf{X}^{(2)}$ and :

- Use $\mathbf{X}^{(1)}$ for clustering,
- Use $\mathbf{X}^{(2)}$ for testing.

Information partitioning approaches

Data fission / data thinning

Idea

Decompose the dataset \mathbf{X} into **two independent datasets** $\mathbf{X}^{(1)}$, $\mathbf{X}^{(2)}$ and :

- Use $\mathbf{X}^{(1)}$ for clustering,
- Use $\mathbf{X}^{(2)}$ for testing.

Known decompositions for Gaussian data $X_1, \dots, X_n \sim \mathcal{N}(\mu, \Sigma)$

Data fission (Leiner *et al.*, 2023)

$$X_i^{(1)} = X_i + \tau W_i$$

$$X_i^{(2)} = X_i + W_i/\tau$$

with $W \sim \mathcal{N}(0, \Sigma)$ and $\tau > 0$.

Data thinning (Neufeld *et al.*, 2024)

$$X_i^{(1)} | X_i = x_i \sim \mathcal{N}(\tau x_i, \tau(1 - \tau)\Sigma)$$

$$X_i^{(2)} = X_i - X_i^{(1)}$$

with $\tau \in (0, 1)$.

Information partitioning approaches

Data fission / data thinning

Idea

Decompose the dataset \mathbf{X} into **two independent datasets** $\mathbf{X}^{(1)}$, $\mathbf{X}^{(2)}$ and :

- Use $\mathbf{X}^{(1)}$ for clustering,
- Use $\mathbf{X}^{(2)}$ for testing.

Known decompositions for Gaussian data $X_1, \text{i.i.d.}, X_n \sim \mathcal{N}(\mu, \Sigma)$

Data fission (Leiner *et al.*, 2023)

$$X_i^{(1)} = X_i + \tau W_i$$

$$X_i^{(2)} = X_i + W_i/\tau$$

with $W \sim \mathcal{N}(0, \Sigma)$ and $\tau > 0$.

Data thinning (Neufeld *et al.*, 2024)

$$X_i^{(1)} | X_i = x_i \sim \mathcal{N}(\tau x_i, \tau(1 - \tau)\Sigma)$$

$$X_i^{(2)} = X_i - X_i^{(1)}$$

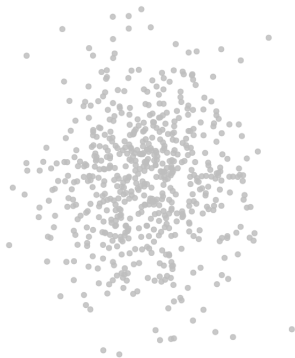
with $\tau \in (0, 1)$.

$$\mathbb{P}_{H_0} \left(\text{Reject } H_0 \text{ based on } \mathbf{X}^{(2)} \mid \text{Select } H_0 \text{ based on } \mathbf{X}^{(1)} \right) = \mathbf{P}_{H_0} \left(\text{Reject } H_0 \text{ based on } \mathbf{X}^{(2)} \right)$$

→ Thanks to **independence**.

Information partitioning approaches

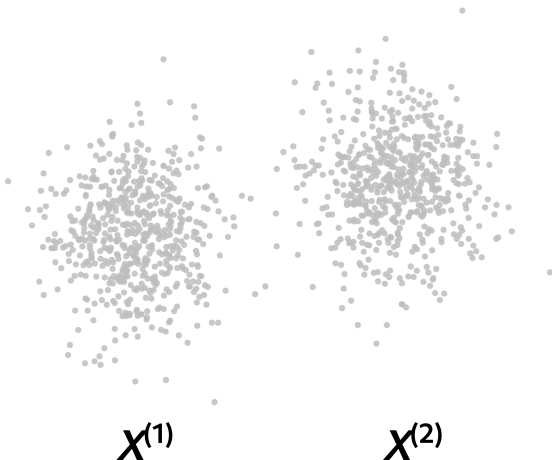
Data fission / data thinning



X

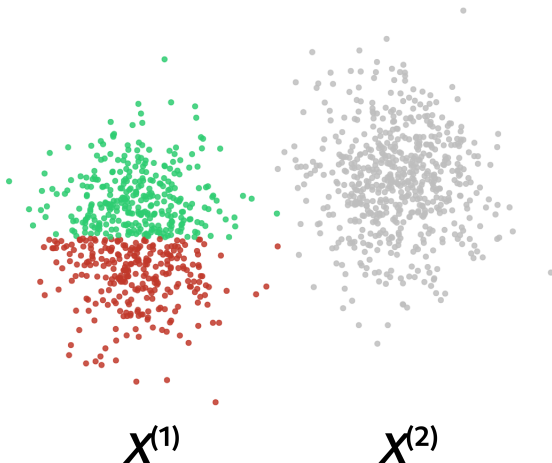
Information partitioning approaches

Data fission / data thinning



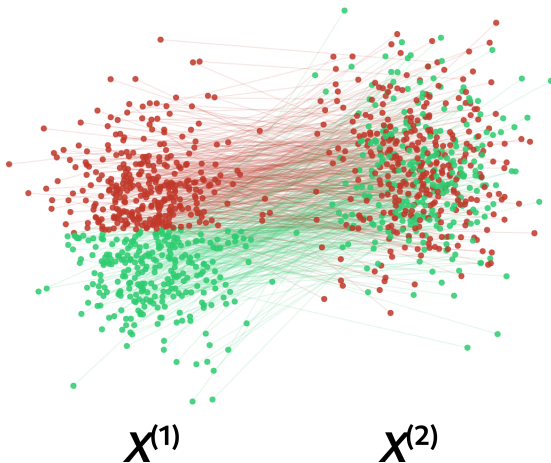
Information partitioning approaches

Data fission / data thinning



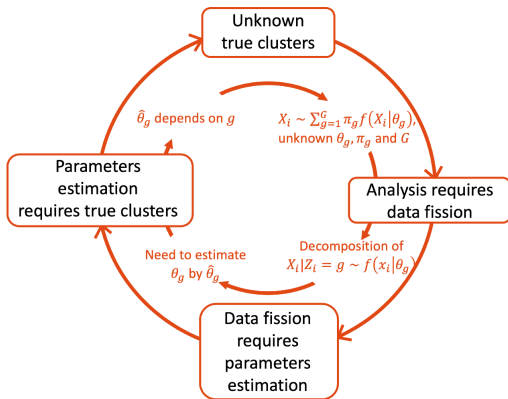
Information partitioning approaches

Data fission / data thinning



Information partitioning approaches

Unsuitable for clustering



Conditional inference

Type I error control

$$\mathbb{P}_{H_0}(\text{Reject } H_0 \text{ based on } \mathbf{X}) \leq \alpha \quad \forall \alpha \in (0, 1).$$

Conditional inference

Selective type I error control (Fithian *et al.* 2014)

$$\mathbb{P}_{H_0}(\text{Reject } H_0 \text{ based on } \mathbf{X} \mid H_0 \text{ selected using } \mathbf{X}) \leq \alpha \quad \forall \alpha \in (0, 1).$$

Conditional inference

Selective type I error control (Fithian *et al.* 2014)

$$\mathbb{P}_{H_0}(\text{Reject } H_0 \text{ based on } \mathbf{X} \mid H_0 \text{ selected using } \mathbf{X}) \leq \alpha \quad \forall \alpha \in (0, 1).$$

Selective type I error control for clustering

$$\mathbb{P}_{H_0}(\text{Reject } H_0 \text{ based on } \mathbf{X} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{X})) \leq \alpha \quad \forall \alpha \in (0, 1),$$

→ $\mathcal{G}_1, \mathcal{G}_2$ are two non-overlapping groups of observations,

→ \mathcal{C} is a clustering algorithm.

Conditional inference

Selective type I error control (Fithian *et al.* 2014)

$$\mathbb{P}_{H_0}(\text{Reject } H_0 \text{ based on } \mathbf{X} \mid H_0 \text{ selected using } \mathbf{X}) \leq \alpha \quad \forall \alpha \in (0, 1).$$

Selective type I error control for clustering

$$\mathbb{P}_{H_0}(\text{Reject } H_0 \text{ based on } \mathbf{X} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{X})) \leq \alpha \quad \forall \alpha \in (0, 1),$$

→ $\mathcal{G}_1, \mathcal{G}_2$ are two non-overlapping groups of observations,

→ \mathcal{C} is a clustering algorithm.

To date, **conditional approaches** are the only methods that have been demonstrated to be suitable for post-clustering inference.

State-of-the-art and goal

State-of-the-art and goal

(Main) state-of-the-art and its model assumptions

- Independent observations $X_i \sim \mathcal{N}_p(\mu_i, \sigma^2 \mathbb{I}_p)$ + estimation of σ .

State-of-the-art and goal

(Main) state-of-the-art and its model assumptions

- Independent observations $X_i \sim \mathcal{N}_p(\mu_i, \sigma^2 \mathbb{I}_p)$ + estimation of σ .
→ Gao *et al.* (2022), Yun and Barber (2023), Chen and Witten (2023).

State-of-the-art and goal

(Main) state-of-the-art and its model assumptions

- Independent observations $X_i \sim \mathcal{N}_p(\mu_i, \sigma^2 \mathbb{I}_p)$ + estimation of σ .
→ Gao *et al.* (2022), Yun and Barber (2023), Chen and Witten (2023).
- General matrix normal model $\mathbf{X} \sim \mathcal{MN}_{n \times p}(\boldsymbol{\mu}, \mathbf{U}, \boldsymbol{\Sigma})$ + estimation of $\boldsymbol{\Sigma}$.
→ González-Delgado *et al.* (2025).

State-of-the-art and goal

(Main) state-of-the-art and its model assumptions

- Independent observations $X_i \sim \mathcal{N}_p(\mu_i, \sigma^2 \mathbb{I}_p)$ + estimation of σ .
→ Gao *et al.* (2022), Yun and Barber (2023), Chen and Witten (2023).
- General matrix normal model $\mathbf{X} \sim \mathcal{MN}_{n \times p}(\boldsymbol{\mu}, \mathbf{U}, \boldsymbol{\Sigma})$ + estimation of $\boldsymbol{\Sigma}$.
→ González-Delgado *et al.* (2025).

For $\mathbf{U} \in \mathbb{R}^{n \times n}$ and $\boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$ positive definite, that means :

- $\text{vec}(\mathbf{X}) \sim \mathcal{N}_{n \times p}(\text{vec}(\boldsymbol{\mu}), \boldsymbol{\Sigma} \otimes \mathbf{U})$,
- $X_i \sim \mathcal{N}_p(\mu_i, U_{ii} \boldsymbol{\Sigma})$ ($\boldsymbol{\Sigma} \leftrightarrow$ dependence between features),
- $X^j \sim \mathcal{N}_n(\mu^j, \boldsymbol{\Sigma}_{jj} \mathbf{U})$ ($\mathbf{U} \leftrightarrow$ dependence between observations).

State-of-the-art and goal

(Main) state-of-the-art and its model assumptions

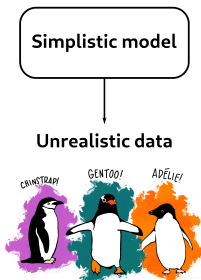
- Independent observations $X_i \sim \mathcal{N}_p(\mu_i, \sigma^2 \mathbb{I}_p)$ + estimation of σ .
→ Gao *et al.* (2022), Yun and Barber (2023), Chen and Witten (2023).
- General matrix normal model $\mathbf{X} \sim \mathcal{MN}_{n \times p}(\boldsymbol{\mu}, \mathbf{U}, \boldsymbol{\Sigma})$ + estimation of $\boldsymbol{\Sigma}$.
→ González-Delgado *et al.* (2025).

Simplistic model

State-of-the-art and goal

(Main) state-of-the-art and its model assumptions

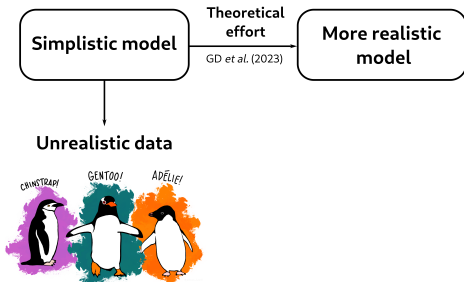
- Independent observations $X_i \sim \mathcal{N}_p(\mu_i, \sigma^2 \mathbb{I}_p)$ + estimation of σ .
→ Gao *et al.* (2022), Yun and Barber (2023), Chen and Witten (2023).
- General matrix normal model $\mathbf{X} \sim \mathcal{MN}_{n \times p}(\boldsymbol{\mu}, \mathbf{U}, \boldsymbol{\Sigma})$ + estimation of $\boldsymbol{\Sigma}$.
→ González-Delgado *et al.* (2025).



State-of-the-art and goal

(Main) state-of-the-art and its model assumptions

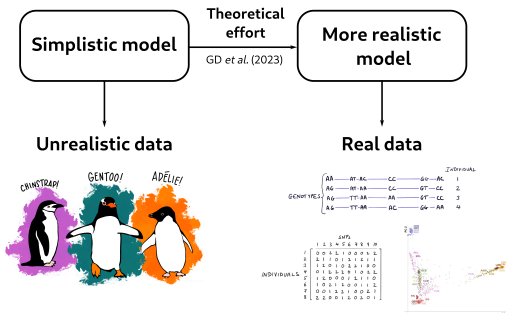
- Independent observations $X_i \sim \mathcal{N}_p(\mu_i, \sigma^2 \mathbb{I}_p)$ + estimation of σ .
→ Gao *et al.* (2022), Yun and Barber (2023), Chen and Witten (2023).
- General matrix normal model $\mathbf{X} \sim \mathcal{MN}_{n \times p}(\boldsymbol{\mu}, \mathbf{U}, \boldsymbol{\Sigma})$ + estimation of $\boldsymbol{\Sigma}$.
→ González-Delgado *et al.* (2025).



State-of-the-art and goal

(Main) state-of-the-art and its model assumptions

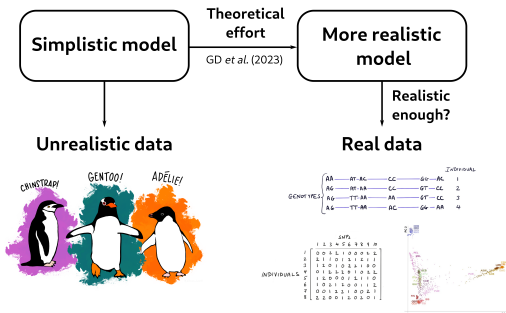
- Independent observations $X_i \sim \mathcal{N}_p(\mu_i, \sigma^2 \mathbb{I}_p) +$ estimation of σ .
→ Gao *et al.* (2022), Yun and Barber (2023), Chen and Witten (2023).
- General matrix normal model $\mathbf{X} \sim \mathcal{MN}_{n \times p}(\boldsymbol{\mu}, \mathbf{U}, \boldsymbol{\Sigma}) +$ estimation of $\boldsymbol{\Sigma}$.
→ González-Delgado *et al.* (2025).



State-of-the-art and goal

(Main) state-of-the-art and its model assumptions

- Independent observations $X_i \sim \mathcal{N}_p(\mu_i, \sigma^2 \mathbb{I}_p) +$ estimation of σ .
→ Gao *et al.* (2022), Yun and Barber (2023), Chen and Witten (2023).
- General matrix normal model $\mathbf{X} \sim \mathcal{MN}_{n \times p}(\boldsymbol{\mu}, \mathbf{U}, \boldsymbol{\Sigma}) +$ estimation of $\boldsymbol{\Sigma}$.
→ González-Delgado *et al.* (2025).



Conditional approaches on real data

Framework setting

Framework setting

- Let $C(\cdot)$ be a clustering algorithm, \mathbf{X} a $n \times p$ random matrix with $\mathbb{E}(\mathbf{X}) = \boldsymbol{\mu}$.

Framework setting

- Let $C(\cdot)$ be a clustering algorithm, \mathbf{X} a $n \times p$ random matrix with $\mathbb{E}(\mathbf{X}) = \boldsymbol{\mu}$.
- Let X_i (resp. μ_i) denote the i -th row of \mathbf{X} (resp. $\boldsymbol{\mu}$) for $i \in [n] = \{1, \dots, n\}$.

Framework setting

- Let $C(\cdot)$ be a clustering algorithm, \mathbf{X} a $n \times p$ random matrix with $\mathbb{E}(\mathbf{X}) = \boldsymbol{\mu}$.
- Let X_i (resp. μ_i) denote the i -th row of \mathbf{X} (resp. $\boldsymbol{\mu}$) for $i \in [n] = \{1, \dots, n\}$.
- For any $\mathcal{G} \subset \{1, \dots, n\}$, let $\bar{X}_{\mathcal{G}} = \frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} X_i$ and $\bar{\mu}_{\mathcal{G}} = \frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} \mu_i$.

Framework setting

- Let $C(\cdot)$ be a clustering algorithm, \mathbf{X} a $n \times p$ random matrix with $\mathbb{E}(\mathbf{X}) = \boldsymbol{\mu}$.
- Let X_i (resp. μ_i) denote the i -th row of \mathbf{X} (resp. $\boldsymbol{\mu}$) for $i \in [n] = \{1, \dots, n\}$.
- For any $\mathcal{G} \subset \{1, \dots, n\}$, let $\bar{X}_{\mathcal{G}} = \frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} X_i$ and $\bar{\mu}_{\mathcal{G}} = \frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} \mu_i$.
- Let $\mathcal{G}_1, \mathcal{G}_2 \subset \{1, \dots, n\}$ be two non-overlapping groups of observations.

Framework setting

- Let $C(\cdot)$ be a clustering algorithm, \mathbf{X} a $n \times p$ random matrix with $\mathbb{E}(\mathbf{X}) = \boldsymbol{\mu}$.
- Let X_i (resp. μ_i) denote the i -th row of \mathbf{X} (resp. $\boldsymbol{\mu}$) for $i \in [n] = \{1, \dots, n\}$.
- For any $\mathcal{G} \subset \{1, \dots, n\}$, let $\bar{X}_{\mathcal{G}} = \frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} X_i$ and $\bar{\mu}_{\mathcal{G}} = \frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} \mu_i$.
- Let $\mathcal{G}_1, \mathcal{G}_2 \subset \{1, \dots, n\}$ be two non-overlapping groups of observations. Considering the column vector $\nu_{\mathcal{G}_1, \mathcal{G}_2} = \nu$ having as components

$$\nu_i = \mathbf{1}\{i \in \mathcal{G}_1\}/|\mathcal{G}_1| - \mathbf{1}\{i \in \mathcal{G}_2\}/|\mathcal{G}_2|,$$

for $i \in [n]$,

Framework setting

- Let $C(\cdot)$ be a clustering algorithm, \mathbf{X} a $n \times p$ random matrix with $\mathbb{E}(\mathbf{X}) = \boldsymbol{\mu}$.
- Let X_i (resp. μ_i) denote the i -th row of \mathbf{X} (resp. $\boldsymbol{\mu}$) for $i \in [n] = \{1, \dots, n\}$.
- For any $\mathcal{G} \subset \{1, \dots, n\}$, let $\bar{X}_{\mathcal{G}} = \frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} X_i$ and $\bar{\mu}_{\mathcal{G}} = \frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} \mu_i$.
- Let $\mathcal{G}_1, \mathcal{G}_2 \subset \{1, \dots, n\}$ be two non-overlapping groups of observations. Considering the column vector $\nu_{\mathcal{G}_1, \mathcal{G}_2} = \nu$ having as components

$$\nu_i = \mathbf{1}\{i \in \mathcal{G}_1\}/|\mathcal{G}_1| - \mathbf{1}\{i \in \mathcal{G}_2\}/|\mathcal{G}_2|,$$

for $i \in [n]$, we can write the difference between the (empirical) group means as

$$\bar{\mu}_{\mathcal{G}_1} - \bar{\mu}_{\mathcal{G}_2} = \boldsymbol{\mu}^T \nu \quad \text{and} \quad \bar{X}_{\mathcal{G}_1} - \bar{X}_{\mathcal{G}_2} = \mathbf{X}^T \nu.$$

Framework setting

- Let $C(\cdot)$ be a clustering algorithm, \mathbf{X} a $n \times p$ random matrix with $\mathbb{E}(\mathbf{X}) = \boldsymbol{\mu}$.
- Let X_i (resp. μ_i) denote the i -th row of \mathbf{X} (resp. $\boldsymbol{\mu}$) for $i \in [n] = \{1, \dots, n\}$.
- For any $\mathcal{G} \subset \{1, \dots, n\}$, let $\bar{X}_{\mathcal{G}} = \frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} X_i$ and $\bar{\mu}_{\mathcal{G}} = \frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} \mu_i$.
- Let $\mathcal{G}_1, \mathcal{G}_2 \subset \{1, \dots, n\}$ be two non-overlapping groups of observations. Considering the column vector $\nu_{\mathcal{G}_1, \mathcal{G}_2} = \nu$ having as components

$$\nu_i = \mathbf{1}\{i \in \mathcal{G}_1\}/|\mathcal{G}_1| - \mathbf{1}\{i \in \mathcal{G}_2\}/|\mathcal{G}_2|,$$

for $i \in [n]$, we can write the difference between the (empirical) group means as

$$\bar{\mu}_{\mathcal{G}_1} - \bar{\mu}_{\mathcal{G}_2} = \boldsymbol{\mu}^T \nu \quad \text{and} \quad \bar{X}_{\mathcal{G}_1} - \bar{X}_{\mathcal{G}_2} = \mathbf{X}^T \nu.$$

We are interested in the following null hypothesis :

$$H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}} : \boldsymbol{\mu}^T \nu = 0. \quad (\text{H0})$$

Standard strategy

1. Choose a model

$$\mathbf{X} \sim \mathcal{MN}_{n \times p}(\boldsymbol{\mu}, \mathbf{U}, \boldsymbol{\Sigma}),$$

where $\boldsymbol{\Sigma}$ can be (over-)estimated and \mathbf{U} is known.

Standard strategy

1. Choose a model

$$\mathbf{X} \sim \mathcal{MN}_{n \times p}(\boldsymbol{\mu}, \mathbf{U}, \boldsymbol{\Sigma}),$$

where $\boldsymbol{\Sigma}$ can be (over-)estimated and \mathbf{U} is known.

→ Under (H0), $\mathbf{X}^T \boldsymbol{\nu} \sim \mathcal{N}_p(\mathbf{0}_p, \mathbf{V})$, with $\mathbf{V} = \boldsymbol{\nu}^T \mathbf{U} \boldsymbol{\nu} \boldsymbol{\Sigma}$.

Standard strategy

1. Choose a model

$$\mathbf{X} \sim \mathcal{MN}_{n \times p}(\boldsymbol{\mu}, \mathbf{U}, \boldsymbol{\Sigma}),$$

where $\boldsymbol{\Sigma}$ can be (over-)estimated and \mathbf{U} is known.

→ Under (H0), $\mathbf{X}^T \boldsymbol{\nu} \sim \mathcal{N}_p(\mathbf{0}_p, \mathbf{V})$, with $\mathbf{V} = \boldsymbol{\nu}^T \mathbf{U} \boldsymbol{\nu} \boldsymbol{\Sigma}$.

→ Therefore, $\|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}}^2 = (\mathbf{X}^T \boldsymbol{\nu})^T \mathbf{V}^{-1} (\mathbf{X}^T \boldsymbol{\nu}) \sim \chi_p^2$ under the null.

Standard strategy

1. Choose a model

$$\mathbf{X} \sim \mathcal{MN}_{n \times p}(\boldsymbol{\mu}, \mathbf{U}, \boldsymbol{\Sigma}),$$

where $\boldsymbol{\Sigma}$ can be (over-)estimated and \mathbf{U} is known.

→ Under (H0), $\mathbf{X}^T \boldsymbol{\nu} \sim \mathcal{N}_p(\mathbf{0}_p, \mathbf{V})$, with $\mathbf{V} = \boldsymbol{\nu}^T \mathbf{U} \boldsymbol{\nu} \boldsymbol{\Sigma}$.

→ Therefore, $\|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}}^2 = (\mathbf{X}^T \boldsymbol{\nu})^T \mathbf{V}^{-1} (\mathbf{X}^T \boldsymbol{\nu}) \sim \chi_p^2$ under the null.

2. Consider a p -value of the form

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}} \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{V}} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{X}), \right)$$

Standard strategy

1. Choose a model

$$\mathbf{X} \sim \mathcal{MN}_{n \times p}(\boldsymbol{\mu}, \mathbf{U}, \boldsymbol{\Sigma}),$$

where $\boldsymbol{\Sigma}$ can be (over-)estimated and \mathbf{U} is known.

→ Under (H0), $\mathbf{X}^T \boldsymbol{\nu} \sim \mathcal{N}_p(\mathbf{0}_p, \mathbf{V})$, with $\mathbf{V} = \boldsymbol{\nu}^T \mathbf{U} \boldsymbol{\nu} \boldsymbol{\Sigma}$.

→ Therefore, $\|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}}^2 = (\mathbf{X}^T \boldsymbol{\nu})^T \mathbf{V}^{-1} (\mathbf{X}^T \boldsymbol{\nu}) \sim \chi_p^2$ under the null.

2. Consider a p -value of the form

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}} \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{V}} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{X}), \right. \\ \left. \boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} \mathbf{X} = \boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} \mathbf{x}, \operatorname{dir}_{\mathbf{V}}(\mathbf{X}^T \boldsymbol{\nu}) = \operatorname{dir}_{\mathbf{V}}(\mathbf{x}^T \boldsymbol{\nu}) \right),$$

where $\boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} = \mathbb{I}_n - \frac{\boldsymbol{\nu} \boldsymbol{\nu}^T}{\|\boldsymbol{\nu}\|_2^2}$ and $\operatorname{dir}_{\mathbf{V}}(u) = u / \|u\|_{\mathbf{V}}$ for all $u \in \mathbb{R}^p$.

Standard strategy

1. Choose a model

$$\mathbf{X} \sim \mathcal{MN}_{n \times p}(\boldsymbol{\mu}, \mathbf{U}, \boldsymbol{\Sigma}),$$

where $\boldsymbol{\Sigma}$ can be (over-)estimated and \mathbf{U} is known.

→ Under (H0), $\mathbf{X}^T \boldsymbol{\nu} \sim \mathcal{N}_p(\mathbf{0}_p, \mathbf{V})$, with $\mathbf{V} = \boldsymbol{\nu}^T \mathbf{U} \boldsymbol{\nu} \boldsymbol{\Sigma}$.

→ Therefore, $\|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}}^2 = (\mathbf{X}^T \boldsymbol{\nu})^T \mathbf{V}^{-1} (\mathbf{X}^T \boldsymbol{\nu}) \sim \chi_p^2$ under the null.

2. Consider a p -value of the form

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}} \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{V}} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{X}), \right. \\ \left. \boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} \mathbf{X} = \boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} \mathbf{x}, \operatorname{dir}_{\mathbf{V}}(\mathbf{X}^T \boldsymbol{\nu}) = \operatorname{dir}_{\mathbf{V}}(\mathbf{x}^T \boldsymbol{\nu}) \right),$$

where $\boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} = \mathbb{I}_n - \frac{\boldsymbol{\nu} \boldsymbol{\nu}^T}{\|\boldsymbol{\nu}\|_2^2}$ and $\operatorname{dir}_{\mathbf{V}}(u) = u / \|u\|_{\mathbf{V}}$ for all $u \in \mathbb{R}^p$.

- The p -value controls the selective type I error for clustering.

Standard strategy

1. Choose a model

$$\mathbf{X} \sim \mathcal{MN}_{n \times p}(\boldsymbol{\mu}, \mathbf{U}, \boldsymbol{\Sigma}),$$

where $\boldsymbol{\Sigma}$ can be (over-)estimated and \mathbf{U} is known.

→ Under (H0), $\mathbf{X}^T \boldsymbol{\nu} \sim \mathcal{N}_p(\mathbf{0}_p, \mathbf{V})$, with $\mathbf{V} = \boldsymbol{\nu}^T \mathbf{U} \boldsymbol{\nu} \boldsymbol{\Sigma}$.

→ Therefore, $\|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}}^2 = (\mathbf{X}^T \boldsymbol{\nu})^T \mathbf{V}^{-1} (\mathbf{X}^T \boldsymbol{\nu}) \sim \chi_p^2$ under the null.

2. Consider a p -value of the form

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}} \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{V}} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{X}), \right. \\ \left. \boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} \mathbf{X} = \boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} \mathbf{x}, \text{dir}_{\mathbf{V}}(\mathbf{X}^T \boldsymbol{\nu}) = \text{dir}_{\mathbf{V}}(\mathbf{x}^T \boldsymbol{\nu}) \right),$$

where $\boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} = \mathbb{I}_n - \frac{\boldsymbol{\nu} \boldsymbol{\nu}^T}{\|\boldsymbol{\nu}\|_2^2}$ and $\text{dir}_{\mathbf{V}}(u) = u / \|u\|_{\mathbf{V}}$ for all $u \in \mathbb{R}^p$.

- The p -value controls the selective type I error for clustering.
- Main challenge : derive an analytically tractable form.

Towards a more interpretable p -value

Rewriting the conditioning set

$$\begin{aligned}\mathbf{X} &= \pi_{\nu}^{\perp} \mathbf{X} + (\mathbb{I}_n - \pi_{\nu}^{\perp}) \mathbf{X} \\ &= \end{aligned}$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$\begin{aligned}\mathbf{X} &= \pi_{\nu}^{\perp} \mathbf{X} + (\mathbb{I}_n - \pi_{\nu}^{\perp}) \mathbf{X} \\ &= \pi_{\nu}^{\perp} \mathbf{X} + \frac{\nu \nu^T}{\|\nu\|_2^2} \mathbf{X}\end{aligned}$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$\begin{aligned}\mathbf{X} &= \pi_{\nu}^{\perp} \mathbf{X} + (\mathbb{I}_n - \pi_{\nu}^{\perp}) \mathbf{X} \\ &= \pi_{\nu}^{\perp} \mathbf{X} + \frac{\nu \nu^T}{\|\nu\|_2^2} \mathbf{X}\end{aligned}$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$\begin{aligned}\mathbf{X} &= \pi_{\nu}^{\perp} \mathbf{X} + (\mathbb{I}_n - \pi_{\nu}^{\perp}) \mathbf{X} \\ &= \pi_{\nu}^{\perp} \mathbf{X} + \frac{\nu}{\|\nu\|_2^2} \nu^T \mathbf{X}\end{aligned}$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$\begin{aligned}\mathbf{X} &= \pi_{\nu}^{\perp} \mathbf{X} + (\mathbb{I}_n - \pi_{\nu}^{\perp}) \mathbf{X} \\ &= \pi_{\nu}^{\perp} \mathbf{X} + \frac{\nu}{\|\nu\|_2^2} \nu^T \mathbf{X} \frac{\|\nu^T \mathbf{X}\|_{\mathbf{V}}}{\|\nu^T \mathbf{X}\|_{\mathbf{V}}}\end{aligned}$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$\begin{aligned}\mathbf{X} &= \pi_{\nu}^{\perp} \mathbf{X} + (\mathbb{I}_n - \pi_{\nu}^{\perp}) \mathbf{X} \\ &= \pi_{\nu}^{\perp} \mathbf{X} + \frac{\nu}{\|\nu\|_2} \nu^T \mathbf{X} \frac{\|\nu^T \mathbf{X}\|_{\mathbf{V}}}{\|\nu^T \mathbf{X}\|_{\mathbf{V}}}\end{aligned}$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$\begin{aligned}\mathbf{X} &= \pi_{\nu}^{\perp} \mathbf{X} + (\mathbb{I}_n - \pi_{\nu}^{\perp}) \mathbf{X} \\ &= \pi_{\nu}^{\perp} \mathbf{X} + \frac{\nu}{\|\nu\|_2} \|\nu^T \mathbf{X}\|_{\mathbf{v}} \text{dir}_{\mathbf{v}}(\nu^T \mathbf{X})\end{aligned}$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$\begin{aligned}\mathbf{X} &= \pi_{\nu}^{\perp} \mathbf{X} + (\mathbb{I}_n - \pi_{\nu}^{\perp}) \mathbf{X} \\ &= \pi_{\nu}^{\perp} \mathbf{X} + \frac{\nu}{\|\nu\|_2^2} \|\nu^T \mathbf{X}\|_{\mathbf{V}} \text{dir}_{\mathbf{V}}(\nu^T \mathbf{X})\end{aligned}$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$\begin{aligned}\mathbf{X} &= \pi_{\nu}^{\perp} \mathbf{X} + (\mathbb{I}_n - \pi_{\nu}^{\perp}) \mathbf{X} \\ &= \pi_{\nu}^{\perp} \mathbf{X} + \frac{\nu}{\|\nu\|_2^2} \|\nu^T \mathbf{X}\|_{\mathbf{V}} \text{dir}_{\mathbf{V}}(\nu^T \mathbf{X})\end{aligned}$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$\begin{aligned}\mathbf{X} &= \pi_\nu^\perp \mathbf{X} + (\mathbb{I}_n - \pi_\nu^\perp) \mathbf{X} \\ &= \pi_\nu^\perp \mathbf{X} + \frac{\nu}{\|\nu\|_2} \|\nu^T \mathbf{X}\|_{\mathbf{V}} \text{dir}_{\mathbf{V}}(\nu^T \mathbf{X})\end{aligned}$$

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\|\mathbf{X}^T \nu\|_{\mathbf{V}} \geq \|\mathbf{x}^T \nu\|_{\mathbf{V}} \mid$$

$$\mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{X}), \pi_\nu^\perp \mathbf{X} = \pi_\nu^\perp \mathbf{x}, \text{dir}_{\mathbf{V}}(\mathbf{X}^T \nu) = \text{dir}_{\mathbf{V}}(\mathbf{x}^T \nu) \right)$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$\begin{aligned}\mathbf{X} &= \boldsymbol{\pi}_\nu^\perp \mathbf{X} + (\mathbb{I}_n - \boldsymbol{\pi}_\nu^\perp) \mathbf{X} \\ &= \boldsymbol{\pi}_\nu^\perp \mathbf{X} + \frac{\boldsymbol{\nu}}{\|\boldsymbol{\nu}\|_2^2} \|\boldsymbol{\nu}^T \mathbf{X}\|_{\mathbf{V}} \text{dir}_{\mathbf{V}}(\boldsymbol{\nu}^T \mathbf{X})\end{aligned}$$

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}} \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{V}} \mid$$

$$\mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C} \left(\boldsymbol{\pi}_\nu^\perp \mathbf{X} + \frac{\boldsymbol{\nu}}{\|\boldsymbol{\nu}\|_2^2} \|\boldsymbol{\nu}^T \mathbf{X}\|_{\mathbf{V}} \text{dir}_{\mathbf{V}}(\boldsymbol{\nu}^T \mathbf{X}) \right), \boldsymbol{\pi}_\nu^\perp \mathbf{X} = \boldsymbol{\pi}_\nu^\perp \mathbf{x}, \text{dir}_{\mathbf{V}}(\mathbf{X}^T \boldsymbol{\nu}) = \text{dir}_{\mathbf{V}}(\mathbf{x}^T \boldsymbol{\nu}) \right)$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}} \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{V}} \mid$$

$$\mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C} \left(\boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} \mathbf{X} + \frac{\boldsymbol{\nu}}{\|\boldsymbol{\nu}\|_2^2} \|\boldsymbol{\nu}^T \mathbf{X}\|_{\mathbf{V}} \text{dir}_{\mathbf{V}}(\boldsymbol{\nu}^T \mathbf{X}) \right), \boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} \mathbf{X} = \boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} \mathbf{x}, \text{dir}_{\mathbf{V}}(\mathbf{X}^T \boldsymbol{\nu}) = \text{dir}_{\mathbf{V}}(\mathbf{x}^T \boldsymbol{\nu}) \right)$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}} \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{V}} \mid \right. \\ \left. \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C} \left(\boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} \mathbf{x} + \frac{\boldsymbol{\nu}}{\|\boldsymbol{\nu}\|_2^2} \|\boldsymbol{\nu}^T \mathbf{X}\|_{\mathbf{V}} \operatorname{dir}_{\mathbf{V}}(\boldsymbol{\nu}^T \mathbf{x}) \right), \boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} \mathbf{X} = \boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} \mathbf{x}, \operatorname{dir}_{\mathbf{V}}(\mathbf{X}^T \boldsymbol{\nu}) = \operatorname{dir}_{\mathbf{V}}(\mathbf{x}^T \boldsymbol{\nu}) \right)$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}} \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{V}} \mid \right. \\ \left. \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C} \left(\boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} \mathbf{x} + \frac{\boldsymbol{\nu}}{\|\boldsymbol{\nu}\|_2^2} \|\boldsymbol{\nu}^T \mathbf{X}\|_{\mathbf{V}} \operatorname{dir}_{\mathbf{V}}(\boldsymbol{\nu}^T \mathbf{x}) \right), \boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} \mathbf{X} = \boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} \mathbf{x}, \underbrace{\operatorname{dir}_{\mathbf{V}}(\mathbf{X}^T \boldsymbol{\nu})}_{\perp \|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}}} = \operatorname{dir}_{\mathbf{V}}(\mathbf{x}^T \boldsymbol{\nu}) \right)$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}} \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{V}} \mid \right.$$

$$\left. \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C} \left(\underbrace{\pi_{\boldsymbol{\nu}}^\perp \mathbf{x}}_{\perp \|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}} \Leftrightarrow \mathbf{U} \in \mathcal{CS}(n)} + \frac{\boldsymbol{\nu}}{\|\boldsymbol{\nu}\|_2^2} \|\boldsymbol{\nu}^T \mathbf{X}\|_{\mathbf{V}} \operatorname{dir}_{\mathbf{V}}(\boldsymbol{\nu}^T \mathbf{x}) \right), \underbrace{\pi_{\boldsymbol{\nu}}^\perp \mathbf{X}}_{\perp \|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}}} = \pi_{\boldsymbol{\nu}}^\perp \mathbf{x}, \underbrace{\operatorname{dir}_{\mathbf{V}}(\mathbf{X}^T \boldsymbol{\nu})}_{\perp \|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}}} = \operatorname{dir}_{\mathbf{V}}(\mathbf{x}^T \boldsymbol{\nu}) \right)$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\|\mathbf{X}^T \boldsymbol{\nu}\|_{\mathbf{V}} \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{V}} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C} \left(\boldsymbol{\pi}_{\boldsymbol{\nu}}^{\perp} \mathbf{x} + \frac{\boldsymbol{\nu}}{\|\boldsymbol{\nu}\|_2^2} \|\boldsymbol{\nu}^T \mathbf{X}\|_{\mathbf{V}} \text{dir}_{\mathbf{V}}(\boldsymbol{\nu}^T \mathbf{x}) \right) \right)$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\nu}} \left(\phi \geq \|\mathbf{x}^T \nu\|_{\mathbf{V}} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C} \left(\pi_{\nu}^{\perp} \mathbf{x} + \frac{\nu}{\|\nu\|_2^2} \phi \operatorname{dir}_{\mathbf{V}}(\nu^T \mathbf{x}) \right) \right)$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\phi \geq \|\mathbf{x}^T \nu\|_{\mathbf{v}} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C} \left(\mathbf{x} + \frac{\nu}{\|\nu\|_2^2} (\phi - \|\mathbf{x}^T \nu\|_{\mathbf{v}}) \operatorname{dir}_{\mathbf{v}}(\mathbf{x}^T \nu) \right) \right)$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\phi \geq \|\mathbf{x}^T \nu\|_{\mathbf{v}} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C} \left(\underbrace{\mathbf{x} + \frac{\nu}{\|\nu\|_2^2} (\phi - \|\mathbf{x}^T \nu\|_{\mathbf{v}}) \operatorname{dir}_{\mathbf{v}}(\mathbf{x}^T \nu)}_{\mathbf{x}'(\phi)} \right) \right)$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\phi \geq \|\mathbf{x}^T \nu\|_{\mathbf{v}} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi)) \right)$$

Towards a more interpretable p -value

Rewriting the conditioning set

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\phi \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{V}} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi)) \right), \quad \text{with } \phi \sim \chi_p.$$

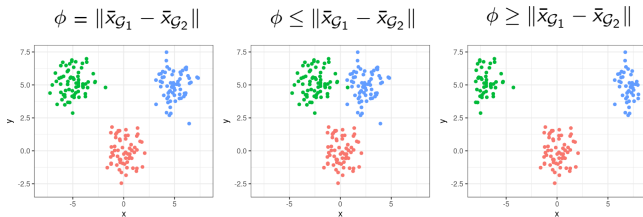
Towards a more interpretable p -value

Rewriting the conditioning set

$$p(\mathbf{x}; \{\mathcal{G}_1, \mathcal{G}_2\}) = \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\phi \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{V}} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi)) \right), \quad \text{with } \phi \sim \chi_p.$$

$\mathbf{x}'(\phi)$ is a perturbation of \mathbf{x} such that

- If $\phi \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{V}}$, \mathcal{G}_1 and \mathcal{G}_2 are pulled apart,
- If $\phi \leq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{V}}$, \mathcal{G}_1 and \mathcal{G}_2 are pushed together.



Interpretation of p -value

A consequence of conditioning

$$\begin{aligned} p(\mathbf{x}) &= \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\phi \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{V}} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi)) \right) \\ &= \frac{\mathbb{E}[\mathbf{1}\{\phi \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{V}}, \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi))\}]}{\mathbb{E}[\mathbf{1}\{\mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi))\}]} \end{aligned}$$

Interpretation of p -value

A consequence of conditioning

$$\begin{aligned} p(\mathbf{x}) &= \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\phi \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{v}} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi)) \right) \\ &= \frac{\mathbb{E}[\mathbf{1}\{\phi \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{v}}, \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi))\}]}{\mathbb{E}[\mathbf{1}\{\mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi))\}]} \end{aligned}$$

Informally

Do I have room to push the clusters together while preserving the clustering?

Interpretation of p -value

A consequence of conditioning

$$\begin{aligned} p(x) &= \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\phi \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{v}} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi)) \right) \\ &= \frac{\mathbb{E}[\mathbf{1}\{\phi \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{v}}, \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi))\}]}{\mathbb{E}[\mathbf{1}\{\mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi))\}]} \end{aligned}$$

Informally

Do I have room to push the clusters together while preserving the clustering?

- Yes \rightarrow reject (H_0) \rightarrow different clusters.
- No \rightarrow accept (H_0) \rightarrow same cluster.

Interpretation of p -value

A consequence of conditioning

$$\begin{aligned} p(x) &= \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\phi \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{v}} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi)) \right) \\ &= \frac{\mathbb{E}[\mathbf{1}\{\phi \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{v}}, \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi))\}]}{\mathbb{E}[\mathbf{1}\{\mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi))\}]} \end{aligned}$$

Informally

Do I have room to push the clusters together while preserving the clustering?

- Yes \rightarrow reject (H_0) \rightarrow different clusters.
- No \rightarrow accept (H_0) \rightarrow same cluster.

In practice

- We need $\{\phi : \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi))\} \neq \emptyset$.

Interpretation of p -value

A consequence of conditioning

$$\begin{aligned} p(x) &= \mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\phi \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{v}} \mid \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi)) \right) \\ &= \frac{\mathbb{E}[\mathbf{1}\{\phi \geq \|\mathbf{x}^T \boldsymbol{\nu}\|_{\mathbf{v}}, \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi))\}]}{\mathbb{E}[\mathbf{1}\{\mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi))\}]} \end{aligned}$$

Informally

Do I have room to push the clusters together while preserving the clustering?

- Yes \rightarrow reject (H_0) \rightarrow different clusters.
- No \rightarrow accept (H_0) \rightarrow same cluster.

In practice

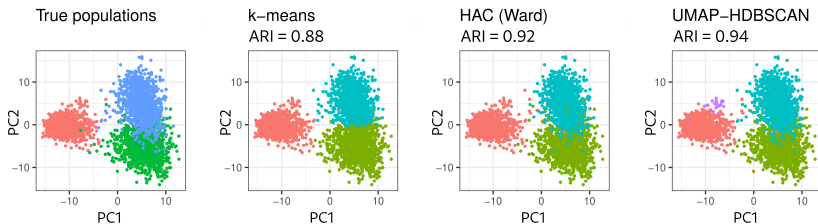
- We need $\{\phi : \mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi))\} \neq \emptyset$.
- If $\{\mathcal{G}_1, \mathcal{G}_2 \in \mathcal{C}(\mathbf{x}'(\phi))\}$ is unlikely, p -values will be very close to 0 or 1.

Consequences of conditioning (I)

Clustering need to be robust to small perturbations

Simulation setting

SNPs for $n = 3000$ individuals equally sampled from three populations¹ + Clustering performed on the first 100 PCs.



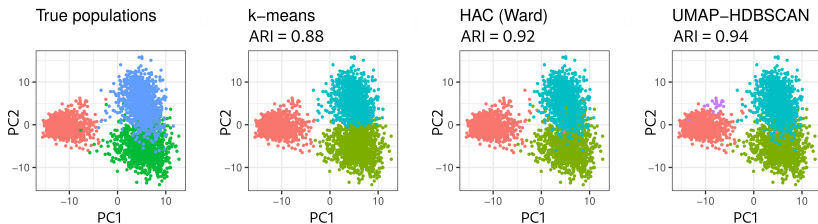
1. Data : Diaz-Papkovich, *unpublished*.

Consequences of conditioning (I)

Clustering need to be robust to small perturbations

Simulation setting

SNPs for $n = 3000$ individuals equally sampled from three populations¹ + Clustering performed on the first 100 PCs.



$\{\mathcal{G}_1, \mathcal{G}_2\}$	$\{\text{red}, \text{green}\}$	$\{\text{red}, \text{cyan}\}$	$\{\text{green}, \text{cyan}\}$
k-means	0.279	0.554	0.178
HAC (Ward)	$< 10^{-16}$	$< 10^{-16}$	$< 10^{-16}$
UMAP+HDBSCAN	1	0.002	0.002

p -values obtained using PCIdep (González-Delgado *et al.* 2025).

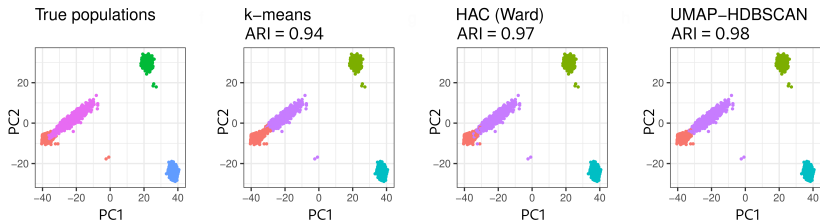
1. Data : Diaz-Papkovich, *unpublished*.

Consequences of conditioning (II)

Shapes and relative positions of clusters may hinder inference

Simulation setting

SNPs for $n = 4000$ individuals equally sampled from four populations : two sources, one unrelated population, and an admixed population derived from the sources¹ + Clustering performed on the first 100 PCs.



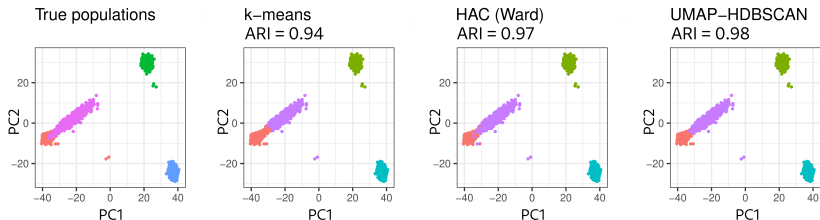
1. Data : Diaz-Papkovich, *unpublished*.

Consequences of conditioning (II)

Shapes and relative positions of clusters may hinder inference

Simulation setting

SNPs for $n = 4000$ individuals equally sampled from four populations : two sources, one unrelated population, and an admixed population derived from the sources¹ + Clustering performed on the first 100 PCs.



$\{\mathcal{G}_1, \mathcal{G}_2\}$	$\{\text{red}, \text{green}\}$	$\{\text{red}, \text{cyan}\}$	$\{\text{green}, \text{cyan}\}$	$\{\text{red}, \text{purple}\}$	$\{\text{green}, \text{purple}\}$	$\{\text{cyan}, \text{purple}\}$
<i>k</i> -means	0.118	0.113	0.005	0.189	0.767	0.667
HAC	0.084	0.011	$< 10^{-16}$	0.343	$< 10^{-16}$	$< 10^{-16}$
UMAP+HDBSCAN	1	1	0.002	0.002	0.002	0.002

p-values obtained using PCIdep (González-Delgado *et al.* 2025).

1. Data : Diaz-Papkovich, *unpublished*.

Conclusion and perspectives

Main message

Conditional inference **inherently limits** post-clustering inference.

Conclusion and perspectives

Main message

Conditional inference **inherently limits** post-clustering inference.

The theory appears unsuitable in multiple practical applications

- Clustering algorithms sensitive to smooth boundaries,
- Realistic cluster placements.

Conclusion and perspectives

Main message

Conditional inference **inherently limits** post-clustering inference.

The theory appears unsuitable in multiple practical applications

- Clustering algorithms sensitive to smooth boundaries,
- Realistic cluster placements.

Limitations come from conditioning

Advances in this theory should proceed through a **new selective inference framework**, such as **simultaneous inference** :

$$\mathbb{P}_{H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}}} \left(\forall \{\mathcal{G}_1, \mathcal{G}_2\} : \text{Reject } H_0^{\{\mathcal{G}_1, \mathcal{G}_2\}} \text{ based on } \mathbf{X} \text{ at level } \alpha \right) \leq \alpha.$$

Ongoing exploration...

(Some) references

Conditional inference methods

- **Independence setting** : L. L. Gao, J. Bien, and D. Witten. Selective inference for hierarchical clustering. *Journal of the American Statistical Association*, 119(545), 332–342, 2024.
- **Extension to feature-level test** : B. Hivert, D. Agniel, R. Thiébaud, and B. P. Hejblum. Post-clustering difference testing : Valid inference and practical considerations with applications to ecological and biological data. *Comput. Statist. Data Anal.*, 193 :107916, May 2024.
- **General matrix normal model** : J. González-Delgado, M. Deronzier, J. Cortés and P. Neuvial. Post-clustering inference under dependence, 2025. <https://arxiv.org/abs/2310.11822>.

Information partitioning approaches

- **Data fission** : J. Leiner, B. Duan, L. Wasserman and A. Ramdas. Data Fission : Splitting a Single Data Point. *Journal of the American Statistical Association*, 120(549), 135–146, 2025.
- **Data thinning** : A. Neufeld, A. Dharamshi, L. L. Gao, and D. Witten. Data thinning for convolution-closed distributions. *Journal of Machine Learning Research*, 25 :1-35, 2024.
- **Limitations in post-clustering inference** : B. Hivert, D. Agniel, R. Thiébaud and B. P. Hejblum. Practical limitations for real-life application of data fission and data thinning in post-clustering differential analysis, 2025. <https://arxiv.org/abs/2405.13591>.

Thank you for your attention !

<https://gonzalez-delgado.github.io/>