

# Gestion des individus manquantes en l'intégration de données omiques: imputation multiple dans le cadre de l'AFM

Valentin Voillet, Laurence Liaubet, Philippe Besse, Magali San Cristobal,  
Ignacio González



Toulouse, 03 décembre 2015

## Intégration de données “omiques”

- étude des tableaux dans lesquels un même ensemble d'individus est décrit par plusieurs groupes de variables

The diagram illustrates a data matrix structure. On the left, the word "individus" is written vertically. The matrix is organized into columns representing different groups of variables: "group 1", "group 2", "...", and "group J". Each group contains a vertical stack of small colored rectangles representing individual data points. The colors of these rectangles vary across groups, with "group 1" and "group 2" showing orange and red, and "group J" showing green. The rows represent individual samples, and the columns represent different groups of variables. The labels "stratum 1", "stratum 2", and "stratum S" are placed between the groups of variables, indicating different strata or layers of data.

## Le but :

- étudier les ressemblances entre des individus en termes de toutes les variables

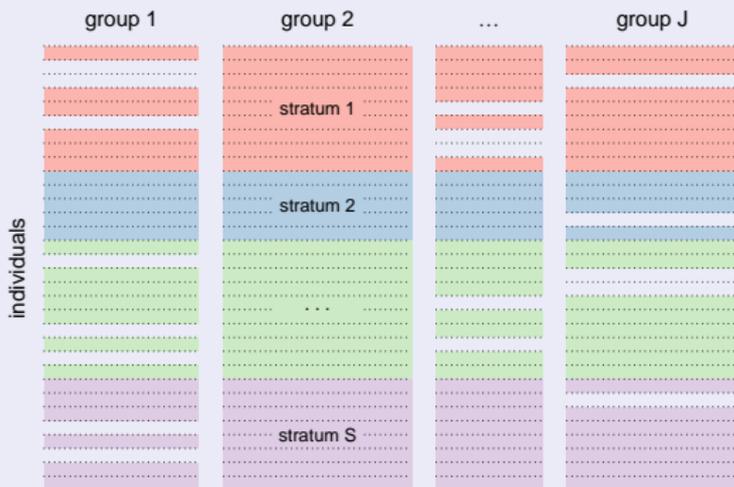
## Problème

- Souvent, on est confronté au problème d'individus manquantes dans un ou plusieurs tableaux



## Problème

- Souvent, on est confronté au problème d'individus manquantes dans un ou plusieurs tableaux



## Problématique des individus manquantes

- la plupart des méthodes statistiques ne peuvent pas être directement appliquées
- dans le cadre de l'analyse de multiple tableaux, peu d'approches sont proposées
- méthodes couramment disponibles pour analyser des données incomplètes :
  - suppression des individus ayant des valeurs manquantes
  - imputation simple, c'est-à-dire les compléter par des valeurs plausibles (par exemple les moyens des cas observés)

## Problèmes des méthodes courantes

- une grande proportion d'individus peut être supprimé
- imputation simple peut déformer les distributions et les relations entre variables
- ces méthodes souffrent de ne pas prendre en compte l'incertitude due aux données manquantes

## Méthodes alternatives

- analyse canonique de corrélations généralisée
- AFM itérative régularisée (RI-MFA)
- imputation multiple en AFM (MI-MFA)

## Problèmes des méthodes courantes

- une grande proportion d'individus peut être supprimé
- imputation simple peut déformer les distributions et les relations entre variables
- ces méthodes souffrent de ne pas prendre en compte l'incertitude due aux données manquantes

## Méthodes alternatives

- analyse canonique de corrélations généralisée
- AFM itérative régularisée (RI-MFA)
- imputation multiple en AFM (MI-MFA)

## Problèmes des méthodes courantes

- une grande proportion d'individus peut être supprimé
- imputation simple peut déformer les distributions et les relations entre variables
- ces méthodes souffrent de ne pas prendre en compte l'incertitude due aux données manquantes

## Méthodes alternatives

- analyse canonique de corrélations généralisée
- AFM itérative régularisée (RI-MFA)
- imputation multiple en AFM (MI-MFA)

## Problèmes des méthodes courantes

- une grande proportion d'individus peut être supprimé
- imputation simple peut déformer les distributions et les relations entre variables
- ces méthodes souffrent de ne pas prendre en compte l'incertitude due aux données manquantes

## Méthodes alternatives

- analyse canonique de corrélations généralisée
- AFM itérative régularisée (RI-MFA)
- **imputation multiple en AFM (MI-MFA)**

## Objectif de la MI-MFA

- estimer les coordonnées des individus sur les premiers composants MFA en présence d'individus manquantes dans un ou plusieurs tableaux
- prendre en compte l'incertitude induite par les individus manquantes sur les composants MFA

## Objectif de la MI-MFA

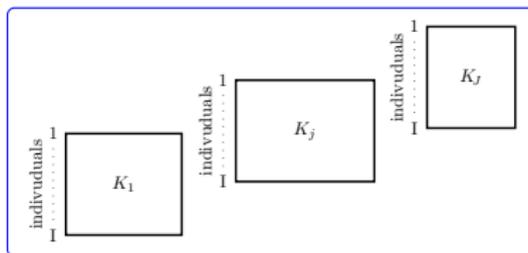
- estimer les coordonnées des individus sur les premiers composants MFA en présence d'individus manquantes dans un ou plusieurs tableaux
- prendre en compte l'incertitude induite par les individus manquantes sur les composants MFA

## Objectif de la MI-MFA

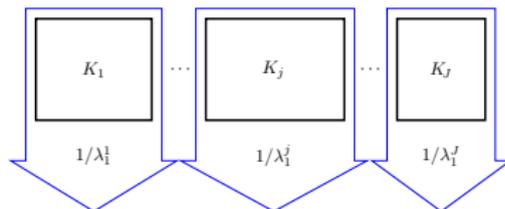
- estimer les coordonnées des individus sur les premiers composants MFA en présence d'individus manquantes dans un ou plusieurs tableaux
- prendre en compte l'incertitude induite par les individus manquantes sur les composants MFA

# L'Analyse Factorielle Multiple (MFA)

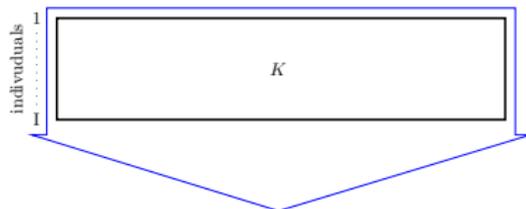
(i)  $J$  data tables



(ii) Separate PCAs and weighting



(iii) Merged and weighted matrix  $K$



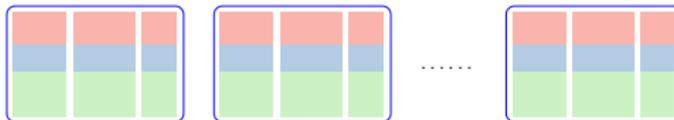
Global PCA

# La méthode MI-MFA

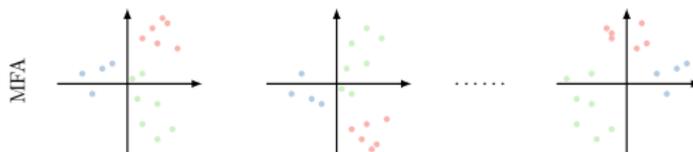
Data with missing rows



$m$  imputed datasets



Analyse each dataset separately

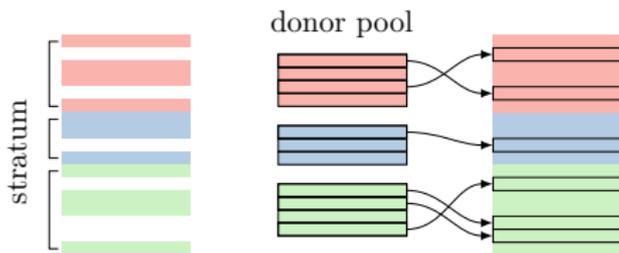


Combine  $m$  configurations into a compromise configuration



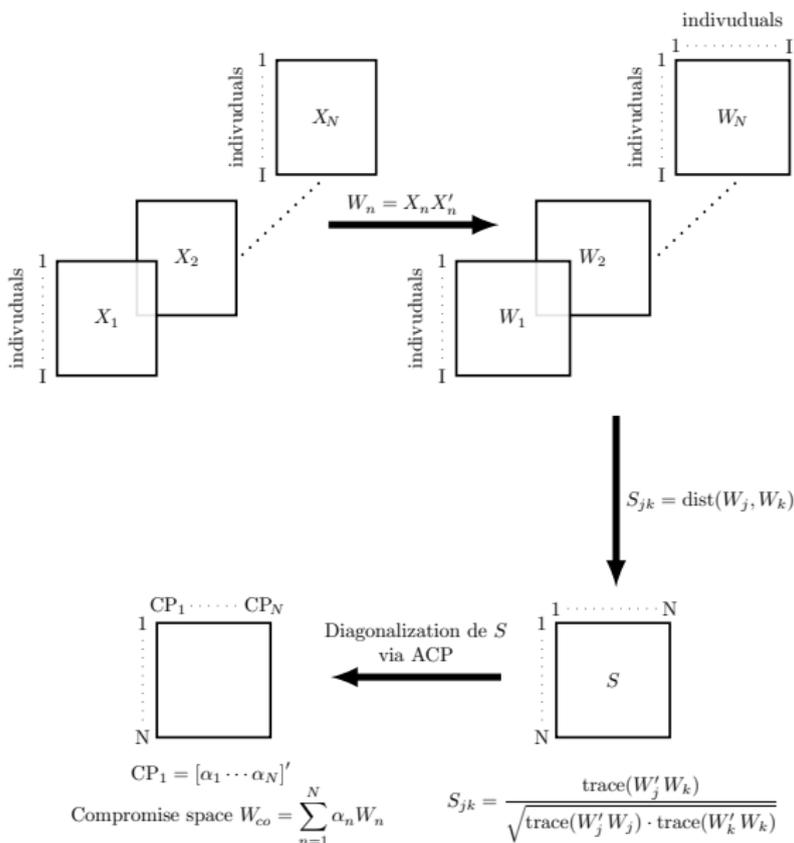
## Comment imputer les individus manquants ?

### La méthode *hot-deck*



- il n'est pas nécessaire de définir un modèle pour la distribution des valeurs manquantes
- prend en compte les ressemblances entre individus et les liaisons entre variables
- peut être appliqué aux données de grands dimensions

# La méthode STATIS (Structuration à Trois Indices de la Statistique)



Elles proviennent d'une étude de toxicité du foie chez la souris  
Les souris ont été exposés à doses de paracétamol dans une expérience contrôlée

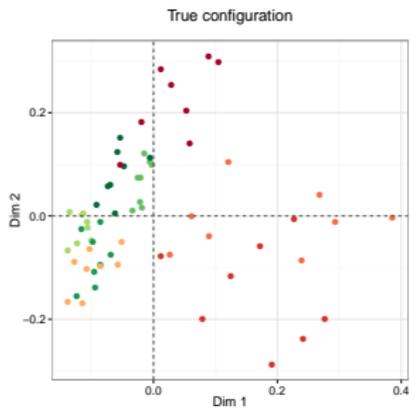
Pour 64 souris, nous disposons :

- des données d'expression de 3116 gènes
- des mesures de 10 variables cliniques contenant marqueurs de lésions du foie
- les 64 souris ont été croisées selon deux facteurs dans un plan à 8 répétitions :
  - la dose : fort (1500 mg/kg ou 2000 mg/kg) et faible (50 mg/kg ou 150 mg/kg)
  - le temps de nécropsie : à 6, 18, 24 et 48 heures

# Les données *liver toxicity*

Condition:

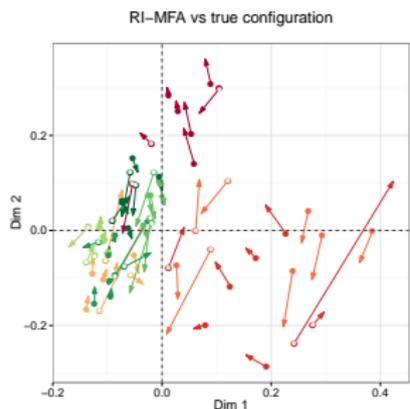
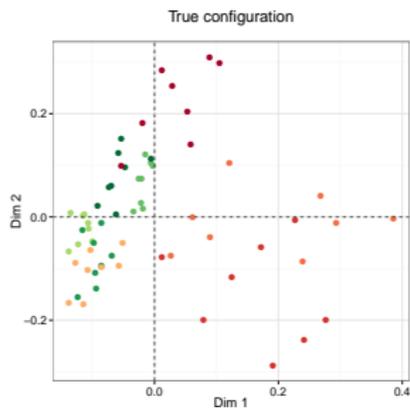
- high-6h
- high-18h
- high-24h
- high-48h
- low-6h
- low-18h
- low-24h
- low-48h



# Les données *liver toxicity*

Condition:

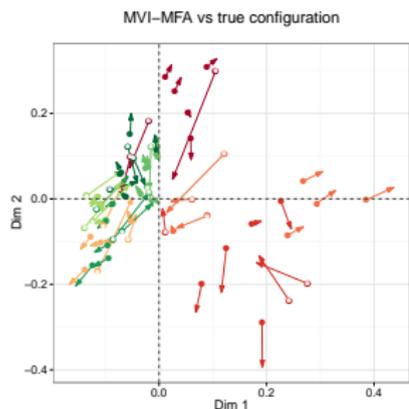
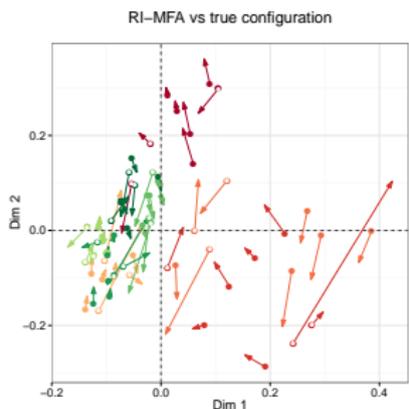
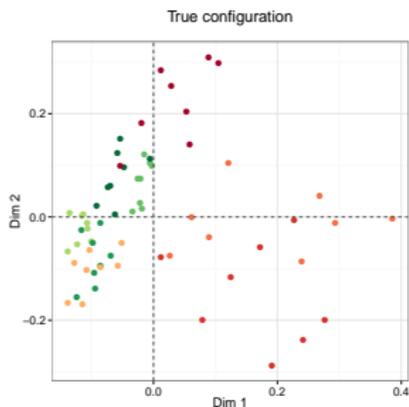
- high-6h
- high-18h
- high-24h
- high-48h
- low-6h
- low-18h
- low-24h
- low-48h



# Les données *liver toxicity*

Condition:

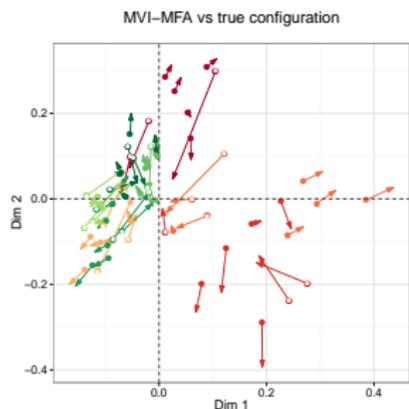
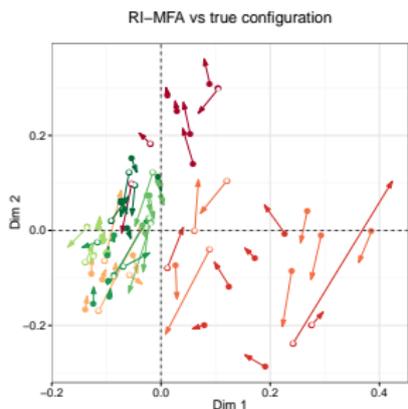
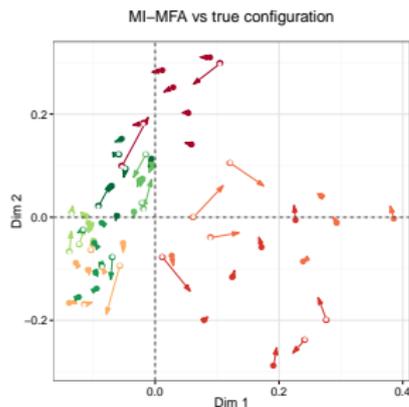
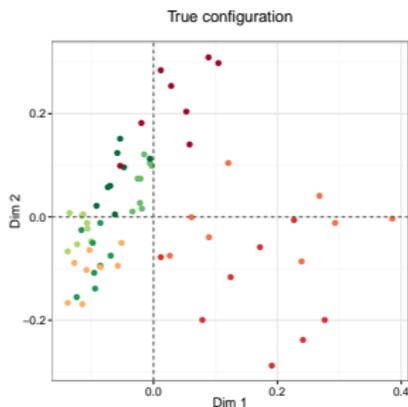
- high-6h
- high-18h
- high-24h
- high-48h
- low-6h
- low-18h
- low-24h
- low-48h



# Les données *liver toxicity*

Condition:

- high-6h
- high-18h
- high-24h
- high-48h
- low-6h
- low-18h
- low-24h
- low-48h



# Les données liverTox

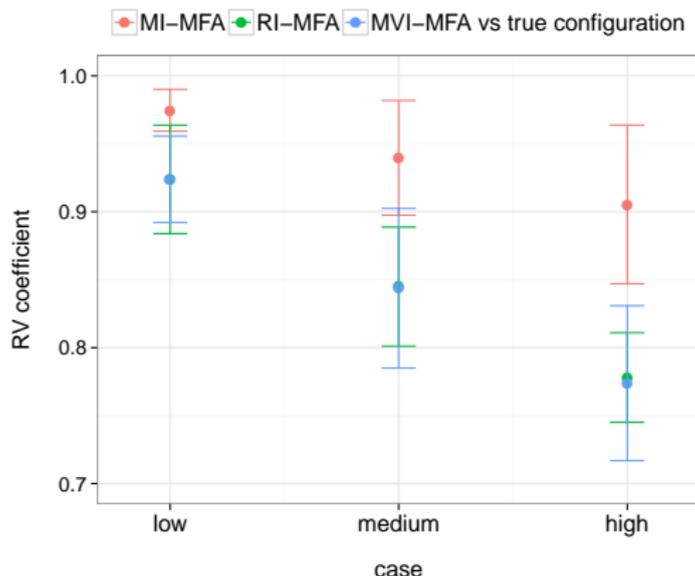
Le nombre d'individus supprimés dans chaque condition sur les données transcriptomiques a été :

low = 1 individu    medium = 2 individus    high = 3 individus

# Les données liverTox

Le nombre d'individus supprimés dans chaque condition sur les données transcriptomiques a été :

low = 1 individu    medium = 2 individus    high = 3 individus



Moyenne des coefficients RV sur 50 configurations 2-dimensionnelles

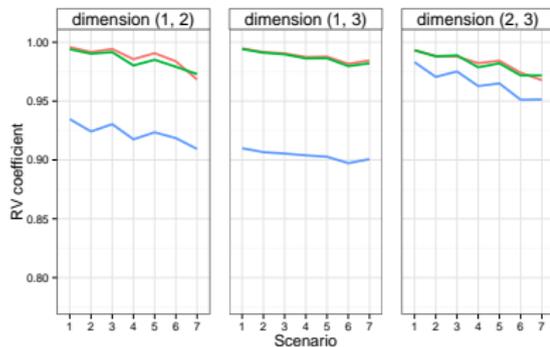
# Les données *liver toxicity*

Différentes scénarios avec des individus manquantes sur les données transcriptomiques et cliniques

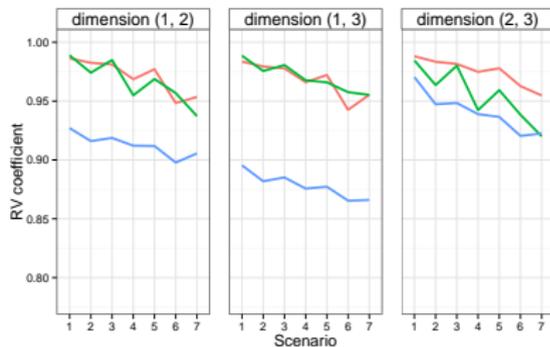
Scénario	Individus manquantes		# cas
	Transcriptomique	Clinique	
1	1	1	56
2	2	1	168
3	1	2	168
4	3	1	280
5	2	2	420
6	3	2	560
7	4	1	280

# Les données liverTox

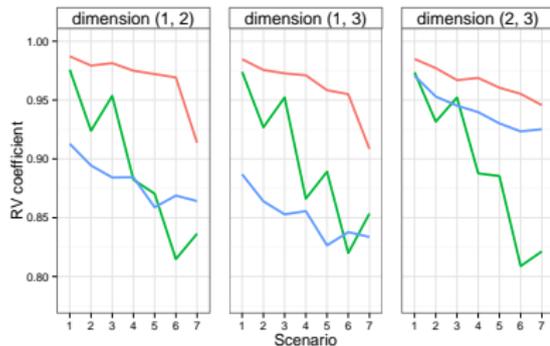
high dose/6h



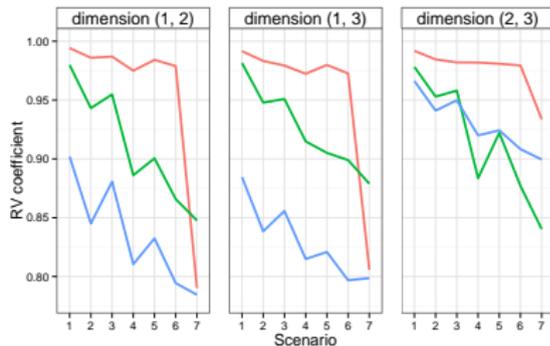
high dose/18h



high dose/24h



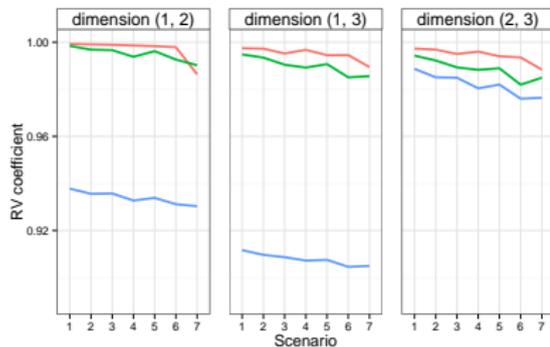
high dose/48h



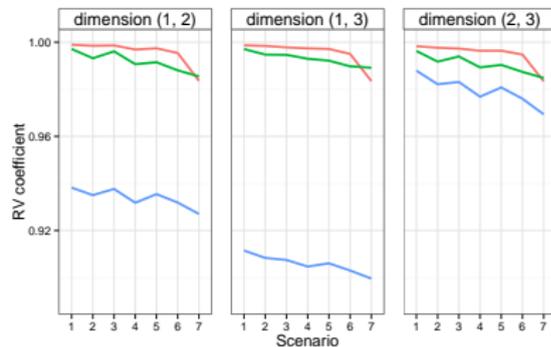
Moyenne des coefficients RV sur 50 configurations

# Les données liverTox

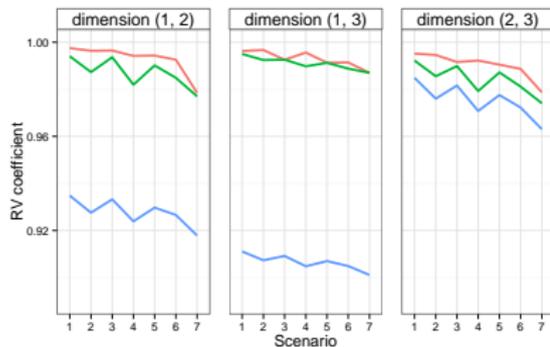
low dose/6h



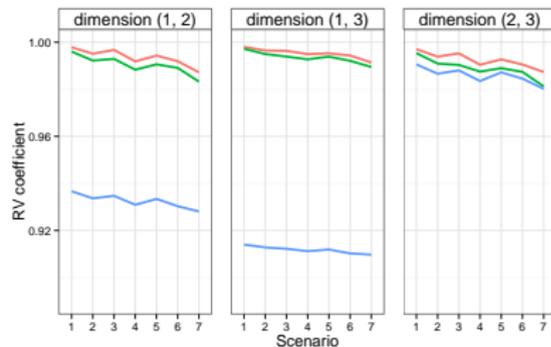
low dose/18h



low dose/24h



low dose/48h



Moyenne des coefficients RV sur 50 configurations

NCI-60 contiennent données transcriptomiques et protéomiques d'une collection de 60 lignes cellulaires cancéreuses :

côlon (7 lignes), rénal (8), ovarien (6), sein (8), prostate (2), poumon (9), système nerveux central (6), leucémies (6) et mélanomes (8)

Pour les 60 lignes cellulaires, nous disposons :

- des données d'expression de 1433 gènes
- des données protéomiques de 162 protéines

# Les données *liver toxicity*

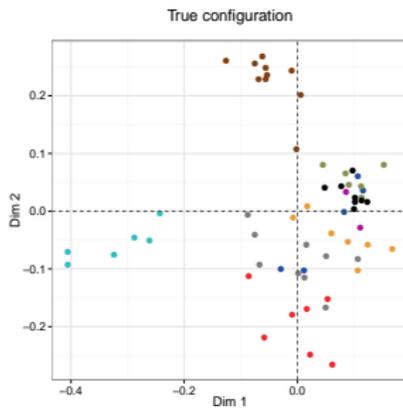
Différentes scénarios avec des individus manquantes sur les données transcriptomiques et protéomiques

Ligne cellulaire	Individus manquantes		# cas
	Transcriptomique	Protéomique	
Sein	1	0	5
SNC	1	1	60
Côlon	2	0	42
Poumon	2	2	3024
Leucémie	1	1	60
Mélanome	2	2	5040
Ovarien	1	1	84
Prostate	0	0	//
Rénal	2	1	504

# Les données *NCI-60*

Type cancer:

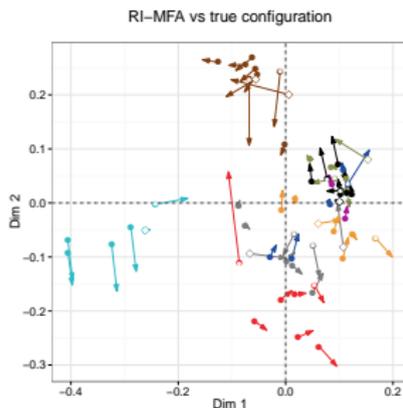
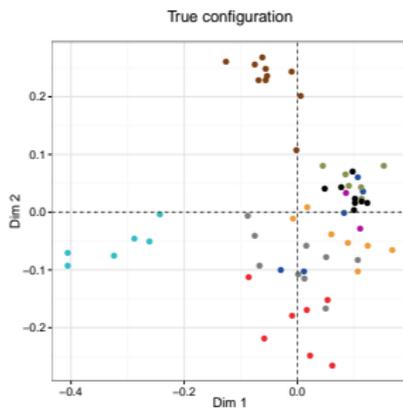
- Colon
- Leucemie
- Melanome
- Ovarien
- Poumon
- Prostate
- Renal
- Sein
- SNC



# Les données *NCI-60*

Type cancer:

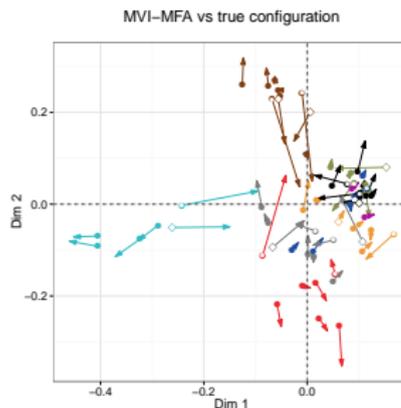
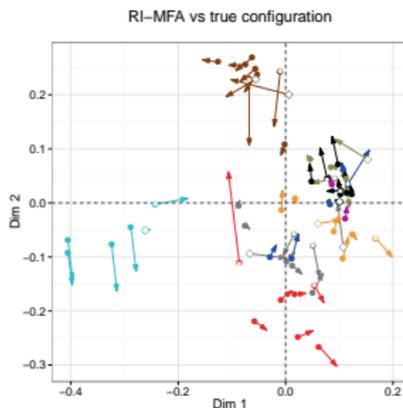
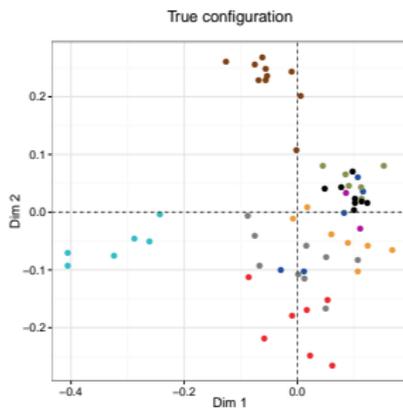
- Colon
- Leucemie
- Melanome
- Ovarien
- Poumon
- Prostate
- Renal
- Sein
- SNC



# Les données NCI-60

Type cancer:

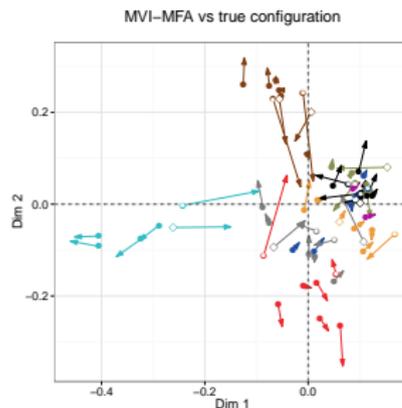
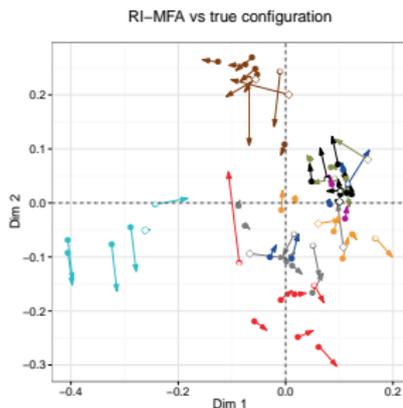
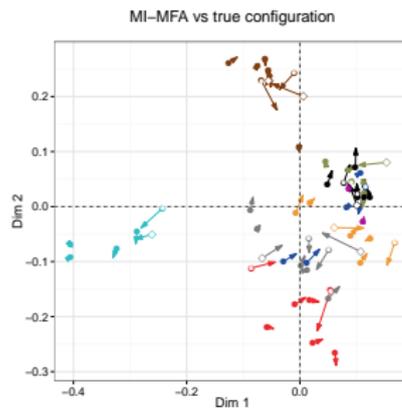
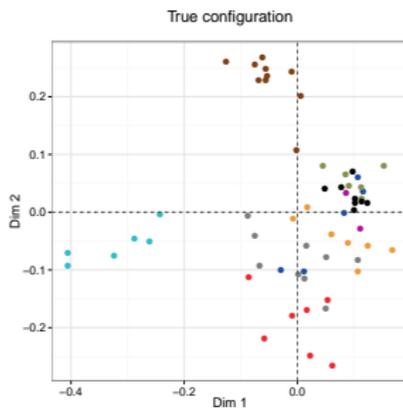
- Côlon
- Leucémie
- Mélanome
- Ovarien
- Poumon
- Prostate
- Rénal
- Sein
- SNC



# Les données *NCI-60*

Type cancer:

- Côlon
- Leucémie
- Mélanome
- Ovarien
- Poumon
- Prostate
- Rénal
- Sein
- SNC



# Pourquoi est-il essentiel l'incertitude ?

Il est toujours possible d'obtenir une configuration d'individus de RI-MFA (ou MVI-MFA) en présence de individus manquantes

Mais, avec seulement cette configuration

- il est impossible savoir si les résultats obtenus sont plausibles
- et si l'utilisateur peut interpréter les résultats correctement
- une configuration unique ne peut pas refléter la variabilité de la estimation

Il est crucial d'avoir un outil permettant d'évaluer la fiabilité des résultats

# Pourquoi est-il essentiel l'incertitude ?

Il est toujours possible d'obtenir une configuration d'individus de RI-MFA (ou MVI-MFA) en présence de individus manquantes

Mais, avec seulement cette configuration

- il est impossible savoir si les résultats obtenus sont plausibles
- et si l'utilisateur peut interpréter les résultats correctement
- une configuration unique ne peut pas refléter la variabilité de la estimation

Il est crucial d'avoir un outil permettant d'évaluer la fiabilité des résultats

# Pourquoi est-il essentiel l'incertitude ?

Il est toujours possible d'obtenir une configuration d'individus de RI-MFA (ou MVI-MFA) en présence de individus manquantes

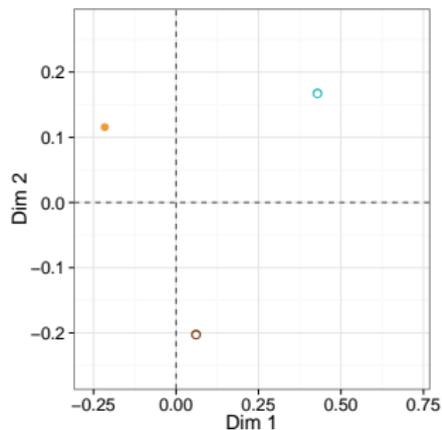
Mais, avec seulement cette configuration

- il est impossible savoir si les résultats obtenus sont plausibles
- et si l'utilisateur peut interpréter les résultats correctement
- une configuration unique ne peut pas refléter la variabilité de la estimation

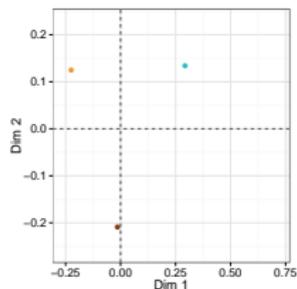
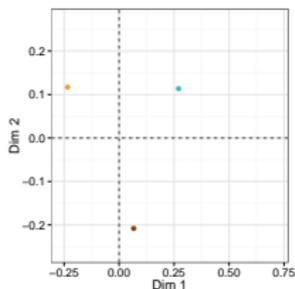
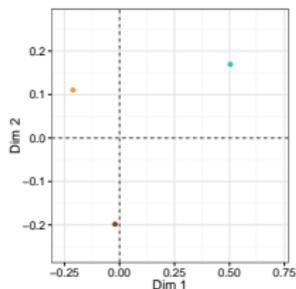
Il est crucial d'avoir un outil permettant d'évaluer la fiabilité des résultats

# Pourquoi est-il essentiel l'incertitude ?

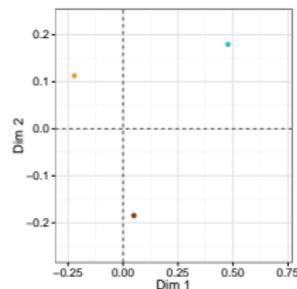
configuration  
compromise



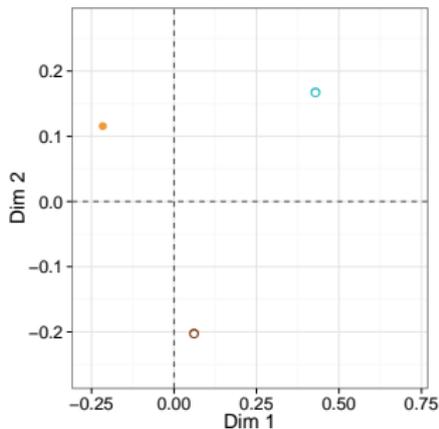
# Pourquoi est-il essentiel l'incertitude ?



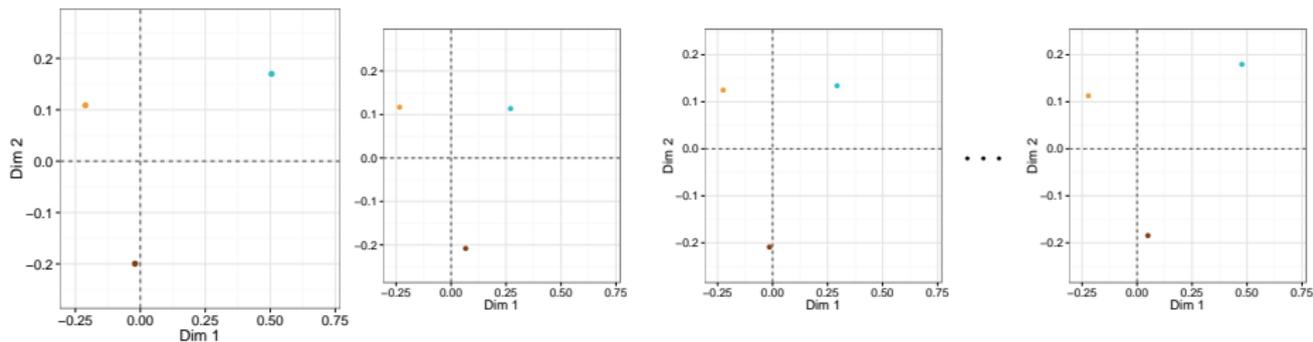
...



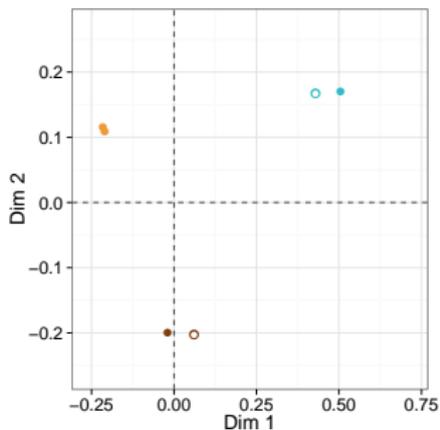
configuration  
compromise



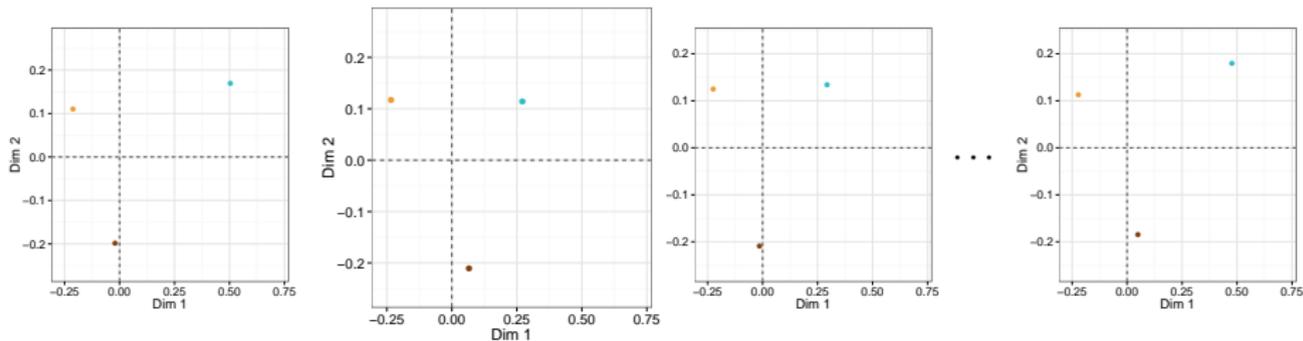
# Pourquoi est-il essentiel l'incertitude ?



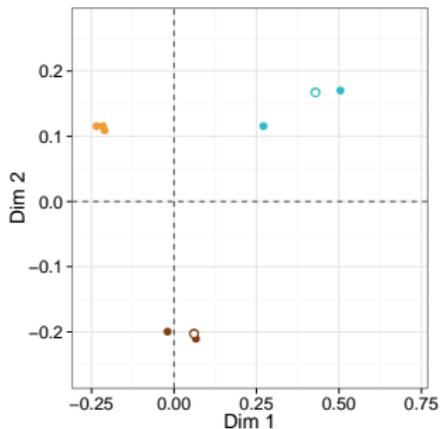
configuration  
compromise



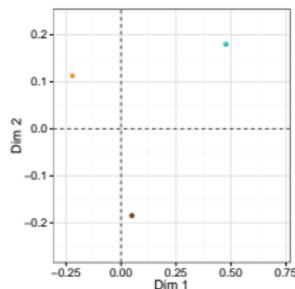
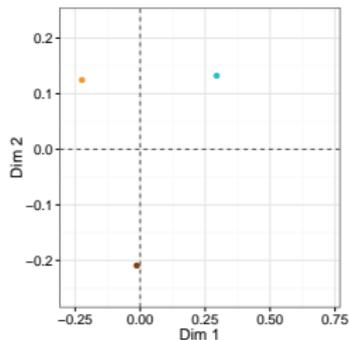
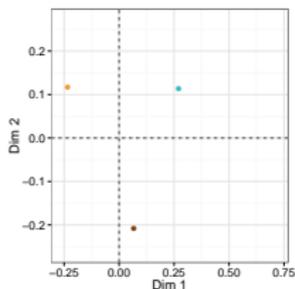
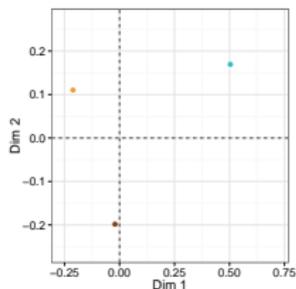
# Pourquoi est-il essentiel l'incertitude ?



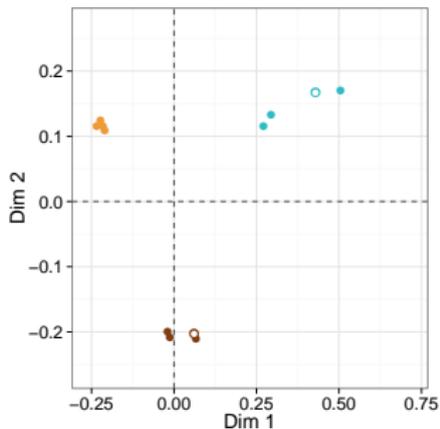
configuration  
compromise



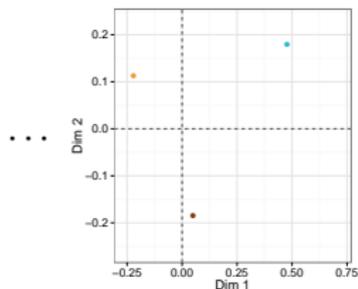
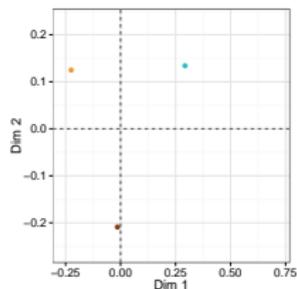
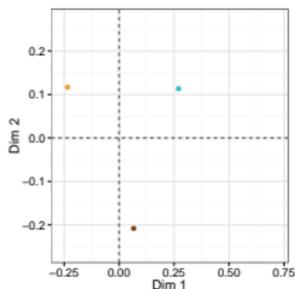
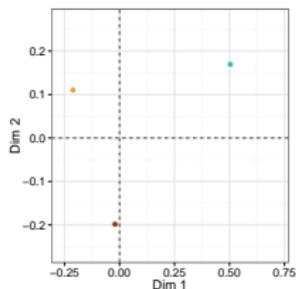
# Pourquoi est-il essentiel l'incertitude ?



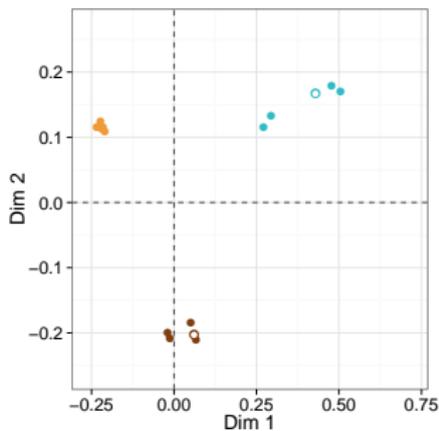
configuration  
compromise



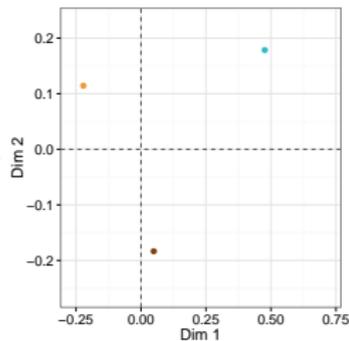
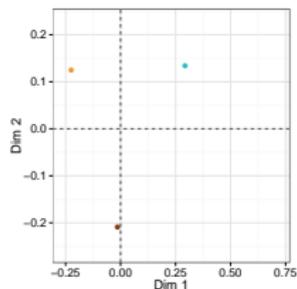
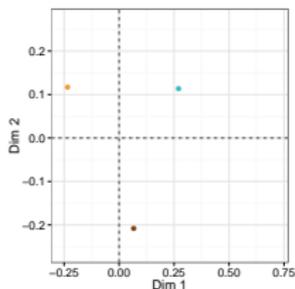
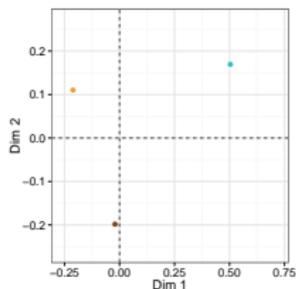
# Pourquoi est-il essentiel l'incertitude ?



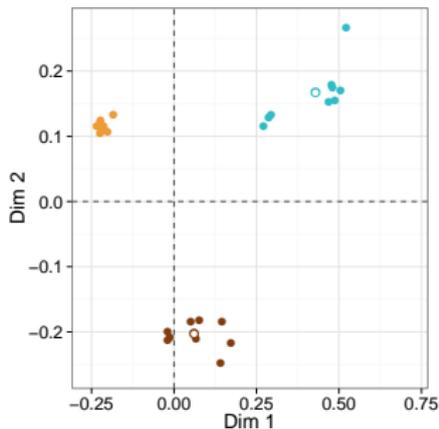
configuration  
compromise



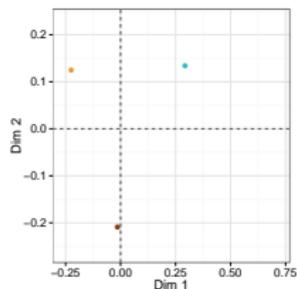
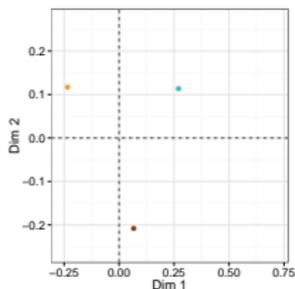
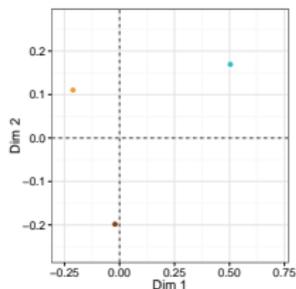
# Pourquoi est-il essentiel l'incertitude ?



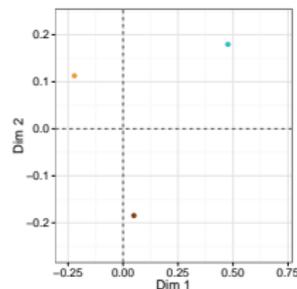
configuration  
compromise



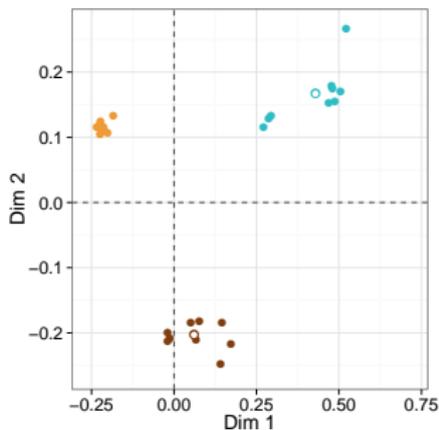
# Pourquoi est-il essentiel l'incertitude ?



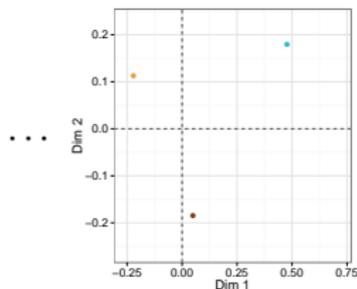
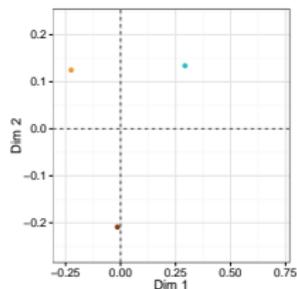
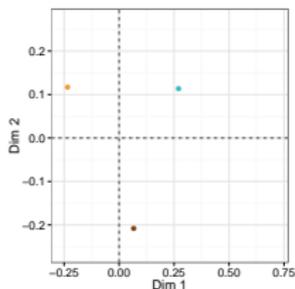
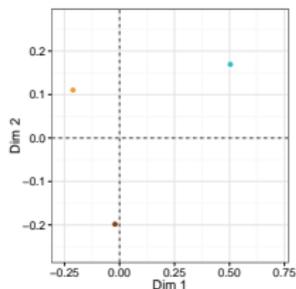
...



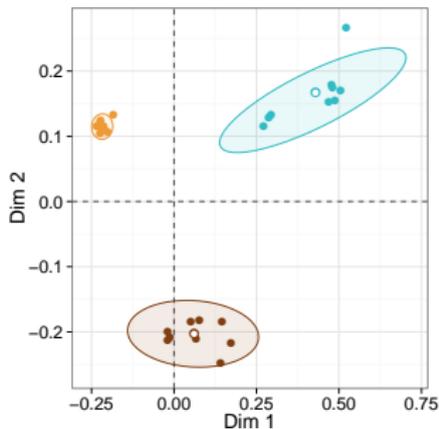
configuration  
compromise



# Pourquoi est-il essentiel l'incertitude ?

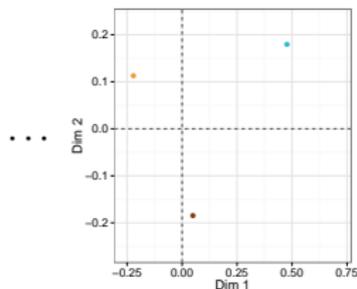
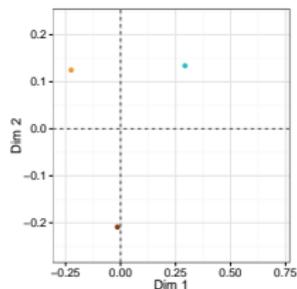
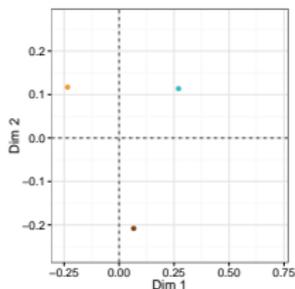
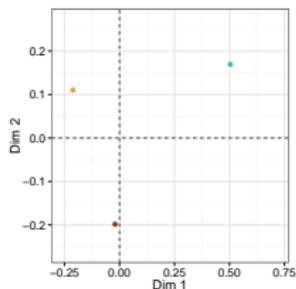


configuration  
compromise

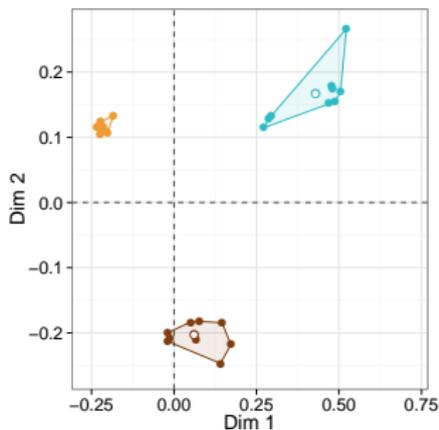


ellipses de confiance

# Pourquoi est-il essentiel l'incertitude ?



configuration  
compromise

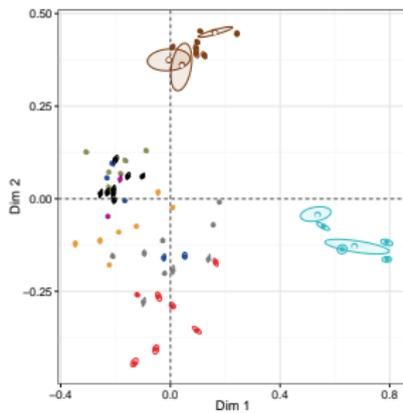


convex hulls

# Pourquoi est-il essentiel l'incertitude ?

Type cancer:

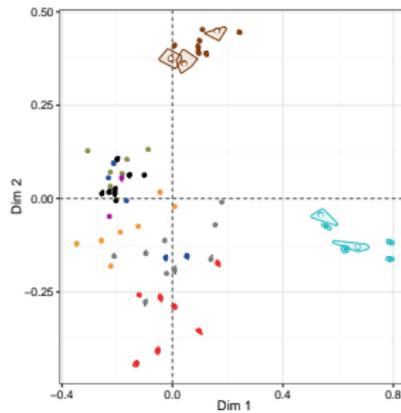
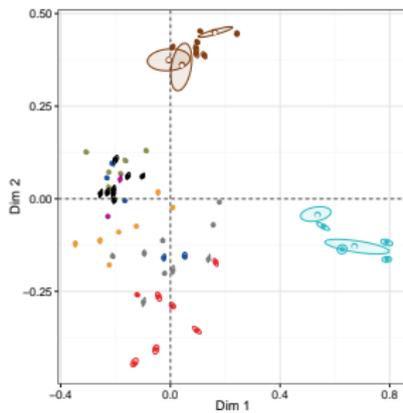
- Colon
- Leucemie
- Melanome
- Ovarien
- Poumon
- Prostate
- Renal
- Sein
- SNC



# Pourquoi est-il essentiel l'incertitude ?

Type cancer:

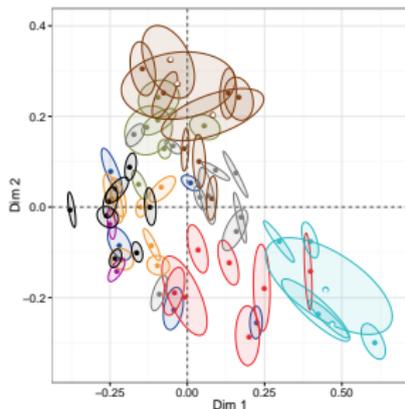
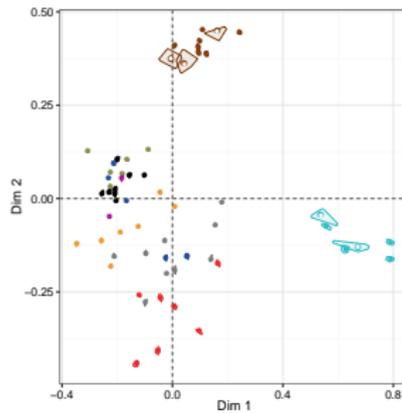
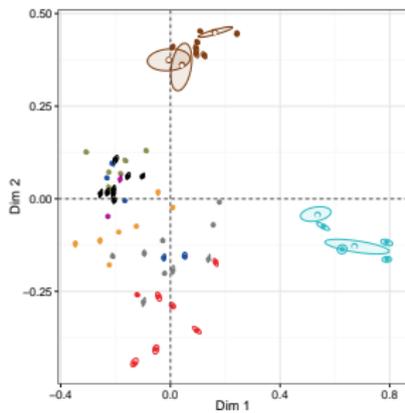
- Colon
- Leucemie
- Melano
- Ovarien
- Poumon
- Prostate
- Renal
- Sein
- SNC



# Pourquoi est-il essentiel l'incertitude ?

Type cancer:

- Colon
- Leucemie
- Melanome
- Ovarien
- Poumon
- Prostate
- Renal
- Sein
- SNC



# Pourquoi est-il essentiel l'incertitude ?

Type cancer:

- Colon
- Leucemie
- Melanome
- Ovarien
- Poumon
- Prostate
- Renal
- Sein
- SNC

