

## Bref récapitulatif de la démarche

Abdoulaye Diabakhaté

**Encadrants:** Stéphane ROBIN et Mahendra MARIADASSOU

19 octobre 2018

# Plan de la présentation

- 1 Introduction
- 2 Les Études
- 3 Piste non approfondie

## Présentation du jeu de donnée utilisé :

*Genocenoses\_env\_parameters\_all\_tara\_initial.xlsx*

Ce jeu de donnée que j'ai reçu de la part de Romain NARCI et sur lequel s'est fait mes études comporte initialement **644 observations sur 20 variables**.

Pour toutes les études menées, on s'est d'abord focalisé sur 15 variables qui sont :

**Lat, Long, T, Sal, chl<sub>a</sub>, O<sub>2</sub>, NO<sub>3m</sub>, NO<sub>3</sub>, NO<sub>2</sub>, NH<sub>4</sub>, SSD, Phos, Si, depth, Fe**

Les **fractions de taille** concernées sont au nombre de 6 :

- 3 fractions composées de 11 matrices de distance : 0-0.2, 0.22-3 et 20-180
- 3 fractions composées de 13 matrices de distance : 5-20, 180-2000 et 0.8-5

## Présentation des matrices de distance :

### Les 11 matrices et les 13 matrices de chaque groupe

Les 11 matrices du premier groupe de fractions sont :

- 1 jaccard\_abundance
- 2 ochiai\_abundance
- 3 sorensen\_abundance
- 4 simka\_jaccard\_abundance
- 5 chord\_hellinger\_prevalence
- 6 jaccard\_canberra\_prevalence
- 7 kulczynski\_prevalence
- 8 ochiai\_prevalence
- 9 whittaker\_prevalence
- 10 simka\_jaccard\_prevalence
- 11 sorensen\_braycurtis\_prevalence

Les 13 matrices du deuxième groupe de fractions sont :

- 1 jaccard\_abundance
- 2 ab\_jaccard\_abundance
- 3 braycurtis\_abundance
- 4 ab\_ochiai\_abundance
- 5 ab\_sorensen\_abundance
- 6 simka\_jaccard\_abundance
- 7 chord\_prevalence
- 8 jaccard\_prevalence
- 9 kulczynski\_prevalence
- 10 ochiai\_prevalence
- 11 whittaker\_prevalence
- 12 simka\_jaccard\_prevalence
- 13 sorensen\_braycurtis\_prevalence

# Syntaxe de la méthode adonis sous R

La syntaxe d'Adonis du package **vegan** de 2 façons est :

- `adonis(matrice ~ variable1 + ... + variable15, data)`
- `adonis(matrice ~ ., data)`

On peut ajouter quelques arguments dans la fonction **adonis** comme :

**permutations=999**, qui est le par défaut, même si on ne le met pas  
ou **method = euclidean**...

## Les analyses menées

Après avoir appliqué notre méthode **adonis**, sur l'ensemble des 15 variables cibles, citées plus haut.

Comme la méthode **adonis** prenait en compte l'ordre d'inclusion des variables, et comme on ne pouvait pas tester tous les ordres d'inclusions possibles ( $2^P - 1$  possibilités), cela nous a motivé à rangé les 15 variables cibles en 3 groupes comme suit :

**Physique**(constituait de **depth** et **T**), **Géographie** (constituait de **Lat** et **Long**)et **Chimie**(constituait de **tout le reste des autres variables**).

C'est à la suite de ce classement qu'on s'est retrouvé à testé **6 ordres de variables des 3 groupes** : ce qui a été fait.

Par la suite on a ajouté une nouvelle variable, cette fois-ci catégorique, **Génocénose** et on s'intéresse à :

- Commencer par tester si l'effet génocénose est significatif une fois les autres variables déjà prise en compte.
- Et, symétriquement, voir quelles variables auraient un effet significatif une fois l'effet génocénose pris en compte.

A la suite de ces analyses on s'est intéressé également à effectuer une **comparaison 2 à 2 des modalités de la variable génocénose**. A cet effet, nous avons utilisé la fonction **pairwise.adonis** du package **pairwiseAdonis**.

Cependant, cette fonction n'admet que la variable facteur, **Génocénose** dans le modèle (en plus de la matrice de distance), et non la totalité de la matrice de design (le reste des covariables), ce qui fera l'objet d'une étude plus approfondie pour les jours qui viennent.



## La fonction `bioenv` (Sélection de variables) : Best Subset Of Environmental Variables With Maximum (Rank) Correlation With Community Dissimilarities

En effet cette fonction, qui s'opère que sur des **variables quantitatives** a pour syntaxe :

La syntaxe de `bioenv` du package **vegan** est :

- `bioenv(matrice ~ variable1 + ... + variable15, data)`
- `bioenv(matrice ~ ., data)`

Cette fonction calcule une matrice de dissimilarité de communauté en utilisant **vegdist**.

Ensuite, elle sélectionne tous les sous-ensembles possibles de variables environnementales, met à l'échelle les variables et calcule les distances euclidiennes pour ce sous-ensemble en utilisant **dist**.

Ensuite, elle trouve la corrélation entre les dissimilarités communautaires et les distances environnementales, et pour chaque taille de sous-ensembles, enregistre le meilleur résultat. Il y'a  $2^P - 1$  sous-ensembles de p-variables.

En effet, un coefficient de corrélation (typiquement le coefficient de **corrélation de rang de Spearman**) est calculé entre les deux matrices et le meilleur sous-ensemble de variables environnementales peut alors être identifié et soumis ensuite à un test de permutation pour déterminer la signification.

La méthode s'arrête lorsque rajouter des variables n'augmente plus la corrélation.