



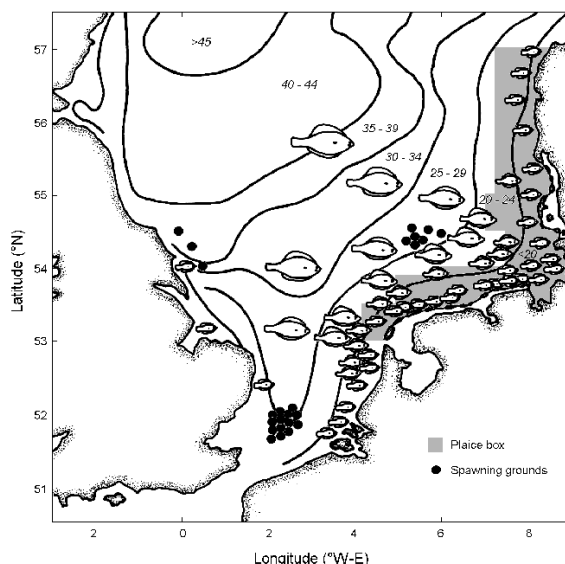
# Formation à R

Modélisation avec les GLM

SÉBASTIEN ROCHETTE, THINKR

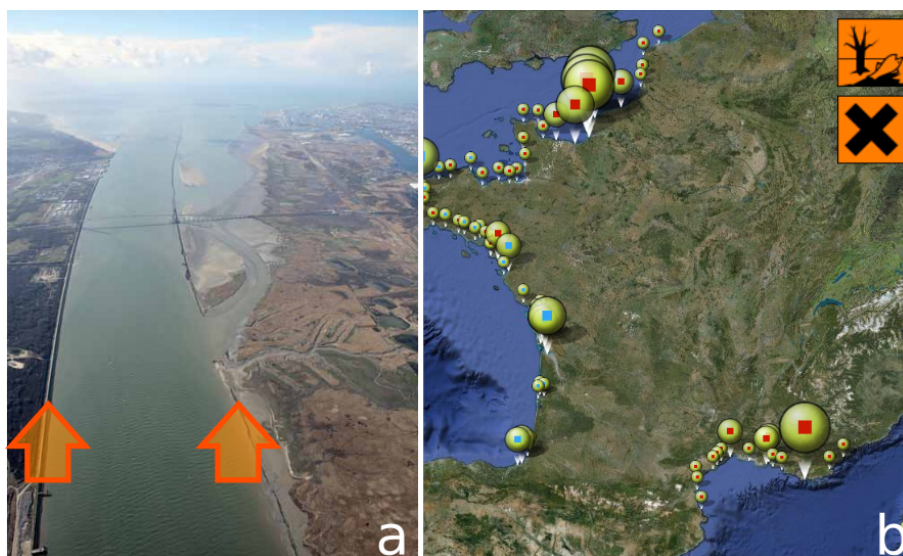
ThinkR

## 1 / 11



**Figure 1** – Plaice box (Rijnsdorp *et al.*)

- Qualité des habitats altérée (Fig. 2b)
- Impact sur le renouvellement des populations
  - Jeune stades = Gouleau d'étranglement
  - La taille et la qualité des nourriceries côtières influent sur la production de juvéniles



**Figure 2** – (a) L'estuaire de la Seine. (b) Niveau de contamination chimique le long des côtes françaises (Ifremer, 2011)

## 2.2. Objectifs

Déterminer les facteurs ayant une influence sur la distribution des poissons plats (*Solea solea*) en Baie de Vilaine et cartographier la distribution moyenne des densités.

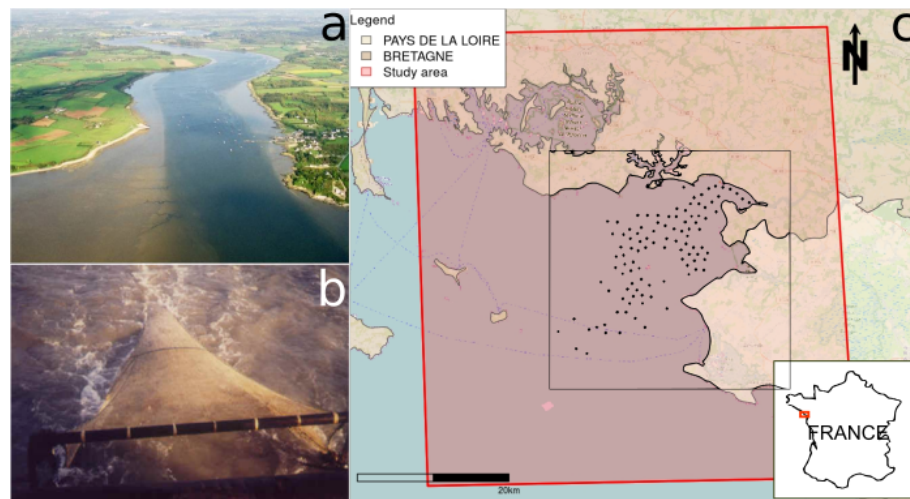
- Cartographier les habitats potentiels nécessite:
  - Connaissance des habitats de juvéniles
  - Campagnes d'échantillonnage dans la zone d'étude
  - Connaissance des covariables environnementales ayant potentiellement de l'influence
    - Cartes exhaustives des covariables environnementales
- Une approche statistique en deux étapes

- Modèle statistique reliant les densités aux covariables
- Prédire les habitats potentiels

## 2.3. Données

Campagne standardisée de chalut à perche dans la baie de Vilaine (Fig. 3)

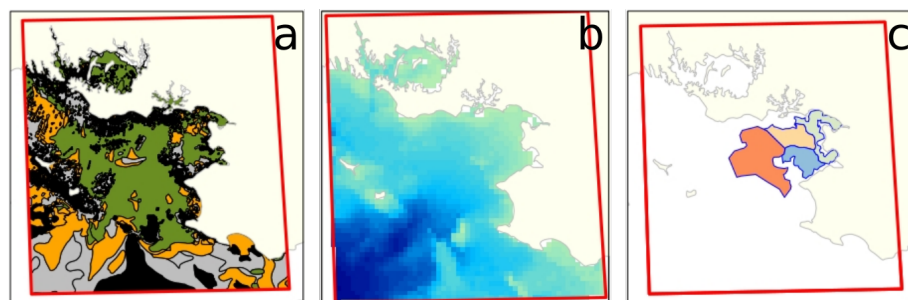
- 1984 – 2010
- En automne
- Juvéniles de l'année (Âge 0)
  - Nb individus / 1000m<sup>2</sup>



**Figure 3** – (a) L'estuaire de la Vilaine. (b) Chalut à perche. (c) Situation des stations d'échantillonnage.

## 2.4. Covariables

- Bathymétrie (Fig. 4a)
  - MNT à 1000m de résolution
  - Projection Mercator
- Structure sédimentaire (Fig. 4b)
  - Fichier shape de polygones
  - Coordonnées géographiques
- Zones biologiques (Fig. 4c)
  - Combinaison bathymétrie, sédiment, habitat
  - Fichier shape de polygones
  - Coordonnées géographiques



**Figure 4** – Covariables en baie de Vilaine. (a) Structure sédimentaire, (b) Bathymétrie et (c) Zones biologiques.

## 2.5. Ajuster un modèle de distribution d'espèces

- Croiser les données avec les cartes de covariables
  - Utiliser un modèle linéaire
- Utiliser les cartes des covariables pour la prédiction (Fig. 5)
  - Une prédiction pour chaque cellule d'un raster

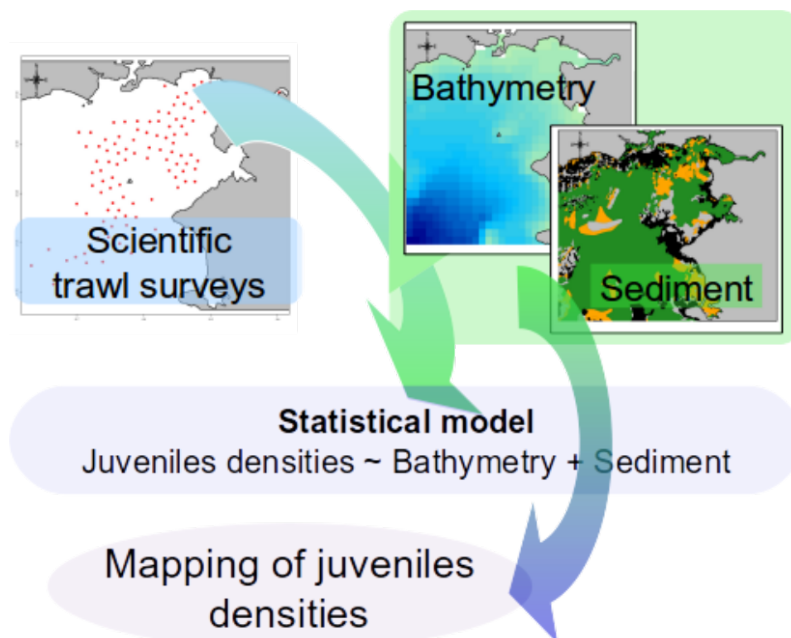


Figure 5 – Procédure pour un modèle de distribution d'espèce

## 2.6. Exploration des données

Prenez le temps d'explorer vos données avant toutes analyses

- Explorer les données et les covariables
  - Explorer le plan d'échantillonnage
  - Explorer les liens potentiels entre les densités et les covariables
  - Explorer les futurs paramètres de modélisation (interactions, distributions)

*Souvenez-vous toujours des objectifs de votre étude !*

*Question : Que recherchons-nous dans cette exploration ?*

## 3. Préparation

### 3.1. Structure des dossiers

*Il convient de toujours conserver les fichiers originaux : les reprojections entraînent toujours quelques pertes, mieux vaut revenir aux originaux lorsque c'est possible.*

L'arborescence de votre dossier de travail est la suivante :

- 01\_Original\_data
  - DEPARTEMENTS
  - Sedim\_GDG\_wgs84



- bathy\_GDG\_1000\_merc (and co)
- Data\_Vilaine\_solea.csv
- 02\_Outputs
- 03\_Figures
- 04\_Functions

### 3.2. Débutons avec R

- Créer un projet Rstudio dans le dossier principal de travail.
- Ouvrez le script R : “Classic\_PresAbs\_Positive\_HSI\_Student.R”
- Lister les différents sous-dossier de travail au début de votre script R

```
# Define working directories -----  
WD <- here()  
# Folder of original files  
origWD <- here("01_Original_data")  
# Folder for outputs  
saveWD <- here("02_Outputs")  
# Folder where to save outputs from R  
figWD <- here("03_Figures")  
# Folder where complementary functions are stored  
funcWD <- here("04_Functions")
```

## 4. Modèle Delta

### 4.1. Étapes

Le modèle sur les données complètes n'était pas satisfaisant. Pour mieux prendre en compte (1) les données d'absences et (2) les fortes valeurs de densités, nous allons utiliser une approche Delta. Le modèle Delta sépare les données en deux sous-groupes, un pour la présence-absence, l'autre pour les densités lorsqu'il y a présence.

- Construction d'un modèle de présence / absence
  - Distribution binomiale
  - Prédiction de probabilités de présence
- Construction d'un modèle sur données positives
  - Distribution à définir
  - Prédiction des densités lorsqu'il y a présence

Puisque les modèles sont ajustés séparément, ils peuvent inclure des covariables différentes.

- L'approche Delta couple les deux sous-modèles
  - Sous-modèle Binomial :  $p_{0/1}$
  - Sous-modèle positif :  $Dens_+$
  - Couplage:  $Density = p_{0/1} \cdot Dens_+$

### 4.2. Sous-modèle sur données positives

#### 4.2.1 Étapes

La procédure à adopter avec le sous-groupe de données est la même qu'avec le jeu de données complet.

- Créer un sous-jeu de données contenant uniquement les observations positives
- Explorer ce nouveau jeu de données

- Explorer les effets potentiels des covariables
- Explorer les potentielles loi de distributions
- Tester les interactions
- Choisir le meilleur mod  le

#### 4.2.2 Validation du mod  le

En utilisant les diff  rents indices pr  sent   pr  c  demment, vous choisissez le mod  le qui s'ajuste le mieux    vos donn  es. C'est donc le meilleur mod  le pour d  crire vos observations. Dans notre cas, nous souhaitons aussi utiliser ce mod  le pour faire de la pr  diction, ce qui n  cessite de s  lectionner un mod  le qui donne de bonnes pr  dictions sur des donn  es non-utilis  es pour l'ajustement du mod  le. Pour cela, nous pouvons utiliser la validation crois  e :

- Ajuster un mod  le sur 90% des donn  es par exemple
- Utiliser le mod  le ajust   pour faire une pr  diction pour les 10% restants
- Comparer les pr  dictions aux observations
- Choisir le mod  le ayant le meilleur indice de comparaison

Vous pourriez utiliser le coefficient de corr  lation entre les pr  dictions et les donn  es de validation comme un indice de qualit   d'ajustement pour s  lectionner le meilleur mod  le. Cependant, l'erreur quadratique moyenne ( $MSE = \text{Mean Squared Error}$ ) voire sa racine ( $RMSE = \text{Root MSE}$ ) est l'indice recommand  . Il mesure la distance moyenne d'une observation    sa pr  diction.

- *Le mod  le s  lectionn   sur la base de l'AIC est-il toujours le meilleur mod  le avec le RMSE ?*

### 4.3. Sous-mod  le Binomial

#### 4.3.1   tapes

La proc  dure    adopter avec le sous-groupe de donn  es est la m  me qu'avec le jeu de donn  es complet.

- Cr  er les observations de pr  sence-absences    partir du jeu de donn  es
- Explorer ce nouveau jeu de donn  es
- Utiliser une distribution binomiale
  - Tester les covariables, les interactions, les fonctions de lien, les crit  res de qualit  
- Choisir le meilleur mod  le

#### 4.3.2 Exploration

#### 4.3.3 Ajuster un mod  le binomial avec une fonction de lien

Le choix de la distribution pour un mod  le de pr  sence-absence est simple, c'est un mod  le binomial. Cependant, un mod  le est g  n  ralement ajust   sur la base de r  sidus Gaussiens. Pour ajuster un mod  le binomial, les donn  es doivent   tre transform  es de telle sorte qu'on puisse ajuster un mod  le lin  aire Gaussien classique dessus. Pour cela, nous utilisons une fonction de lien. La fonction de lien classique d'un mod  le binomial est la fonction logit, mais ce n'est pas la seule. Vous pouvez tester cloglog, probit ou cauchit.

#### 4.3.4 Qualit   d'ajustement d'un mod  le binomial

Une mesure couramment utilis  e pour la qualit   d'ajustement d'un mod  le binomial est "l'aire sous la courbe" (AUC : Area Under the Curve). Un objectif des mod  les binomiaux   tant de pr  dire un succ  s ou un   chec, et non pas seulement une probabilit   de succ  s, on peut vouloir d  finir un seuil (intuitivement 0.5 par exemple) qui transforme la probabilit   de pr  sence en pr  sence ou absence. L'AUC est en quelque sorte une probabilit   de classer correctement les pr  sences et absences. Une d  finition plus compl  te serait :

La probabilité moyenne pour qu'une observation=1 et une observation=0 choisies de manière aléatoire dans le jeu de données montrent une probabilité de présence prédite supérieure pour l'observation=1 par rapport à celle de l'observation=0

Ainsi,  $AUC = 1$  montrerait un modèle "parfait", mais  $AUC = 0.5$  montrerait un modèle plus mauvais que le hasard.

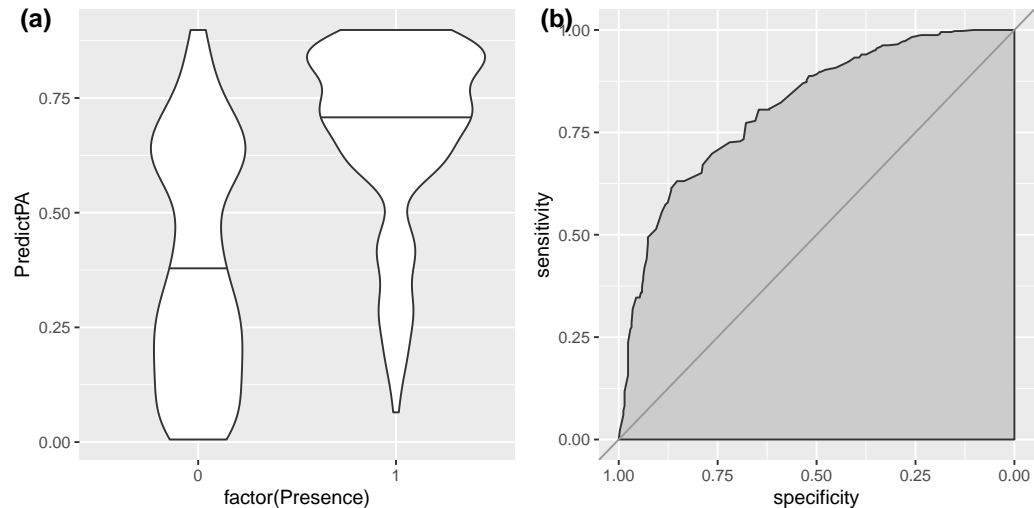


Figure 6 – (a) Prédiction vs Observations. (b) Courbe ROC d'un modèle binomial

#### 4.3.5 Choix du meilleur seuil

- *Le modèle sélectionné sur la base de l'AIC est-il toujours le meilleur modèle avec l'AUC sur les données de validation ?*

### 4.4. Couplage des deux sous-modèles

#### 4.4.1 L'approche Delta

L'approche Delta est la méthode pour coupler les deux sous-modèles. En réalité, les deux modèles sont simplement multipliés l'un à l'autre.

Le couplage des deux sous-modèles c'est :

- Sous-modèle binomial:  $p_{0/1} \sim Bathymetry + Sediment$
- Sous-modèle positif:  $Dens_+ \sim Bathymetry$ 
  - Si log-transformation:  $Dens_+ = \exp(\log(Y_+)) \times \exp(-0.5 \cdot \sigma^2 \cdot \log(Y_+))$
- Couplage:  $Density = p_{0/1} \cdot Dens_+$

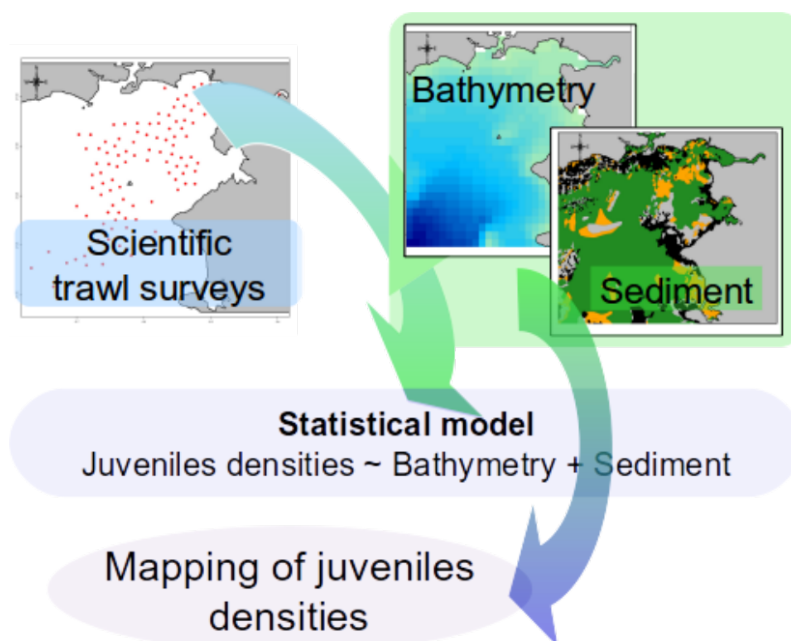
## 5. Modèle d'habitat

### 5.1. Étapes

La réalisation d'une carte de distribution d'espèce (Fig. 7) nécessite :

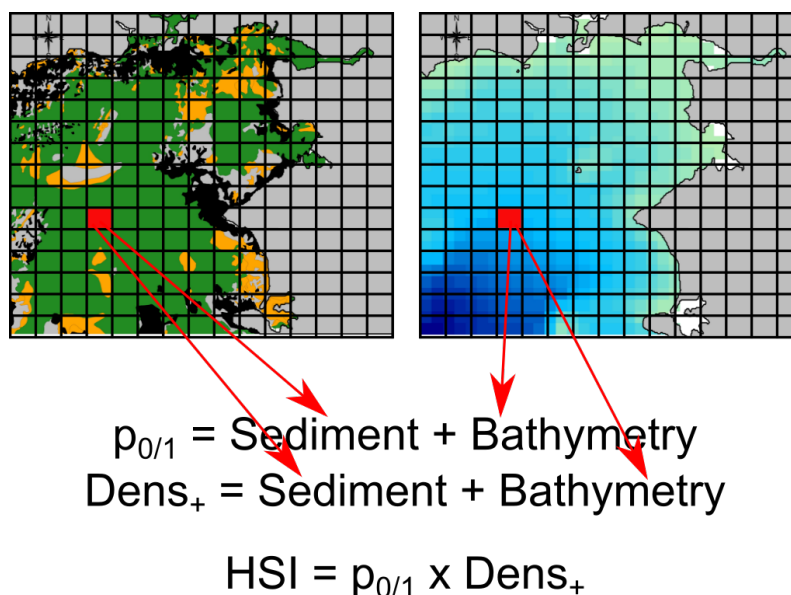
- Un modèle d'habitat potentiel
  - Indice de qualité d'habitat
  - Modèle:  $Density \sim Bathymetry + Sediment$
- Les cartes complètes des covariables
- Une carte des prédictions du modèle





**Figure 7** – Proc  dure pour cartographier une distribution d'esp  ce

Une mani  re simple de r  aliser la carte des pr  dictions est de cr  er un raster qui rassemble l'information de toutes les cartes des covariables n  cessaires, puis d'utiliser le mod  le pour pr  dire dans chaque cellule du raster (Fig. 8).



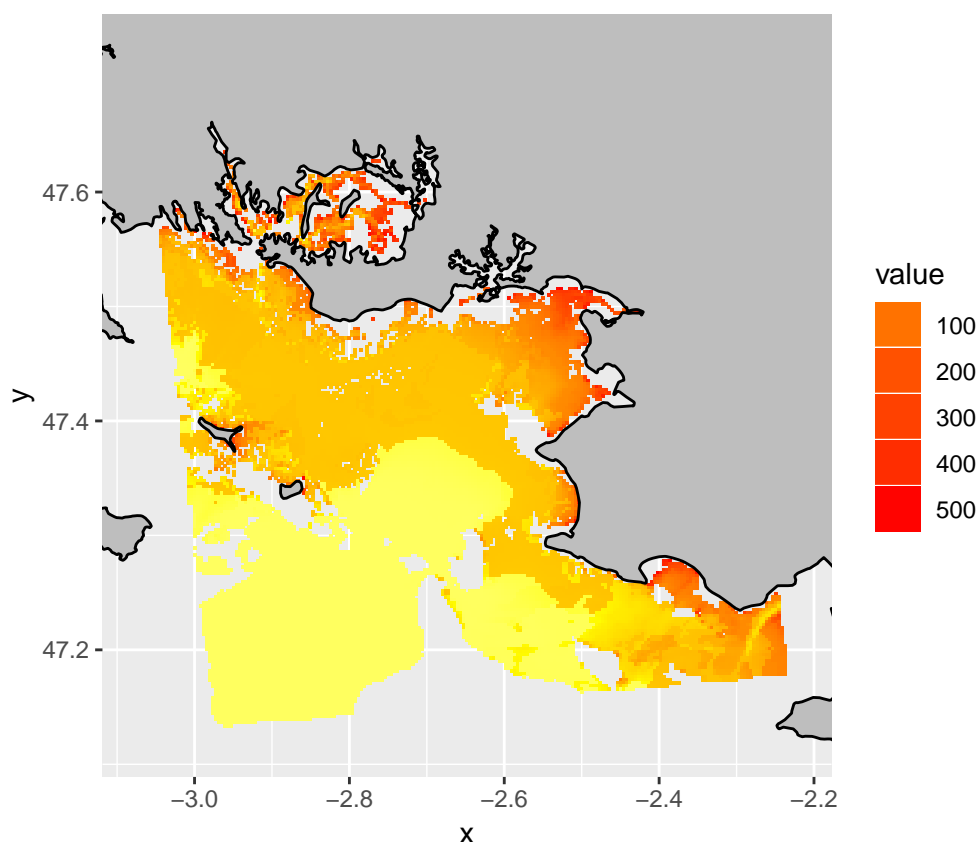
**Figure 8** – Pr  diction de densit   pour chaque cellule d'une carte au format raster

## 5.2. Pr  paration des donn  es

Pour pouvoir faire les pr  dictions dans le raster, ses couches doivent avoir le m  me nom que les colonnes du jeu de donn  es. De plus, un raster est une matrice de valeurs num  riques, ce qui oblige    convertir les covariables en classe en valeurs num  riques, de telle sorte que les niveaux de facteur du raster correspondent    ceux des donn  es, et donc existent dans les mod  les. *Soyez prudents avec la conversion vers des valeurs num  riques, les donn  es doivent rester au format facteur et ne doivent pas   tre utilis  es comme valeurs num  riques dans les mod  les !*

### 5.3. Prédictions

La fonction `predict` a une méthode pour pouvoir être utilisée directement sur un objet Raster (Fig. 9).



**Figure 9** – Prédiction des densités de soles en baie de Vilaine. Échelle de couleur en fonction des quantiles des données originales (10%, 50%, 75%, 95%).

## 6. Conclusion

### 6.1. Modélisation

- **Importance de l'exploration des données**
  - Validation des données
  - Loi de distribution
  - Options pour les modèles
- **Étude des modèles : une approche itérative**
  - Choix de la loi de distribution
  - Choix des combinaisons de covariables
  - Vérification des hypothèses
    - Analyse des résidus
    - Critères de qualité d'ajustement
- **Gardez toujours vos objectifs en tête !**

### 6.2. Modèle Delta

- **Utilité d'un modèle Delta**
  - Les données brutes ne peuvent être modélisées
  - Sens biologique

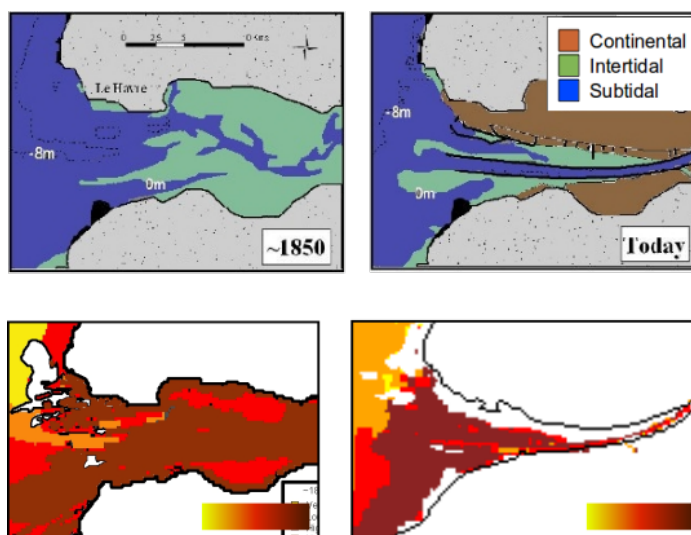
- Présence & densités: pas forcément les mêmes covariables
- **Utilisation d'un modèle Delta**
  - Prédiction utiles
  - Les paramètres des sous-modèles n'ont pas de sens dans le modèle couplé
- **Alternatives**
  - Autres modèles zero-inflated ? Distribution tweedie ?
  - GAM (Attention aux données nécessaires et à l'interprétation)
  - Régression quantile (habitat préférentiels)
  - Random forest (Attention à votre question)

### 6.3. Modèles de distribution d'espèces

- **Outils utiles pour les connaissances biologiques et pour la gestion**
  - Modèle Delta approprié pour les données d'espèces marines
  - Fiable si utilisé avec précaution
    - Ce ne sont que des corrélations...
- **Cette formation est un exemple simple**
  - D'autres perspectives
    - Ajouter une covariable biotique
    - Approche multi-spécifique
    - Pressions anthropiques ?
- **D'autres outils existent**
  - Votre question détermine l'outil à utiliser

### 6.4. Exemples d'applications

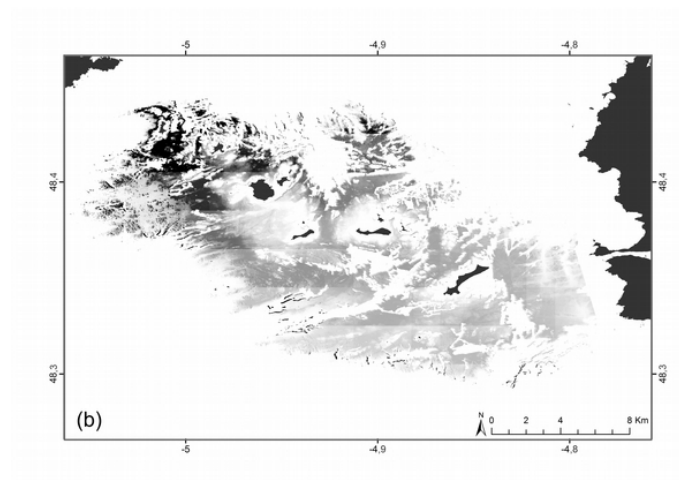
- Effet de la destruction d'habitat sur la biomasse de juvéniles (Fig. 10)
  - Cas de la sole commune dans l'estuaire de Seine
  - Comparaison entre 1850 et 2004
    - Perte en surface : 33%
    - Perte en biomasse : 42%



**Figure 10** – Modélisation des effets de la destruction d'habitats sur la biomasse de sole en Seine. Rochette, S., Rivot, E., Morin, J., Mackinson, S., Riou, P., Le Pape, O. (2010). Effect of nursery habitat degradation on flatfish population renewal. Application to *Solea solea* in the Eastern Channel (Western Europe). Journal of sea Research, 64 : 34-44.

- Estimation de stock pour la gestion (Fig. 11)
  - Cas des laminaires du Parc marin d'Iroise
  - Estimation des biomasses
  - Validation avec les pêcheurs

- Proposition de gestion spatialisée



**Figure 11** – Estimation spatialisée des biomasses de laminaires dans le parc marin d'Iroise pour la gestion de la ressource. Bajjouk T., Rochette S., Ehrhold A., Laurans M., Le Niliot P. (2015). Multi-approach mapping to help spatial planning and management of the kelp species *L. digitata* and *L. hyperborea*: Case study of the Molène archipelago, Brittany. *Journal of Sea Research*.