



# Formation à R

Modélisation avec les GLM

SÉBASTIEN ROCHETTE, THINKR

ThinkR

## Table des mati  res

<b>1</b>	<b>Pr��face</b>	<b>1</b>
<b>2</b>	<b>Pr��sentation de l��tude</b>	<b>1</b>
2.1	Contexte . . . . .	1
2.2	Objectifs . . . . .	2
2.3	Donn��es . . . . .	2
2.4	Covariables . . . . .	3
2.5	Ajuster un mod��le de distribution d��esp��ces . . . . .	3
2.6	Exploration des donn��es . . . . .	4
<b>3</b>	<b>Pr��paration</b>	<b>4</b>
3.1	Structure des dossiers . . . . .	4
3.2	D��butons avec R . . . . .	5
<b>4</b>	<b>Exploration des donn��es</b>	<b>5</b>
4.1	��tapes . . . . .	5
4.2	Liste des diff��rentes ��tapes . . . . .	5
<b>5</b>	<b>Mod��lisation</b>	<b>7</b>
5.1	��tapes . . . . .	7
5.2	Interpr��ter les sorties de mod��les . . . . .	8
5.3	Trouver le meilleur mod��le . . . . .	11
5.4	Pr��dictions du mod��le . . . . .	12

## Pr  face

*La version d'origine de cette formation a   t   cr   e par Olivier Le Pape et   tienne Rivot    Agrocampus Ouest (Rennes, France). Depuis mon doctorat dans leur   quipe, je mets    jour constamment cette formation au gr   de ma recherche et de l  volution du logiciel R.*

*# Generated with R and rmarkdown: Roadmap version - Teacher*

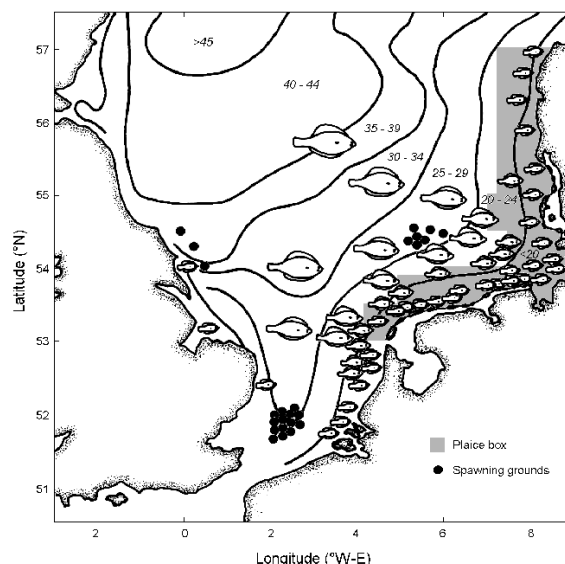
## Pr  sentation de l  tude

*Le contexte et les objectifs de votre   tude d  finissent le type de mod  lisation que vous allez mettre en place sur votre jeu de donn  es.*

Ici, nous utilisons les mod  les lin  aires g  n  ralis  s pour produire une carte de distribution moyenne de la nourricerie de soles communes de la baie de Vilaine.

## Contexte

- Les zones c  ti  res et les estuaires sont des habitats halieutiques essentiels
  - Zones    forte production
  - Nourriceries
  - Zones restreintes avec de fortes densit  s (Fig. 1)
- Pression anthropique   lev  e
  - Perte de surface disponibles (Fig. 2a)
  - Qualit   des habitats alt  r  e (Fig. 2b)
- Impact sur le renouvellement des populations
  - Jeune stades = Goulet d'  tranglement
  - La taille et la qualit   des nourriceries c  ti  res influent sur la production de juv  niles



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



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

## Objectifs

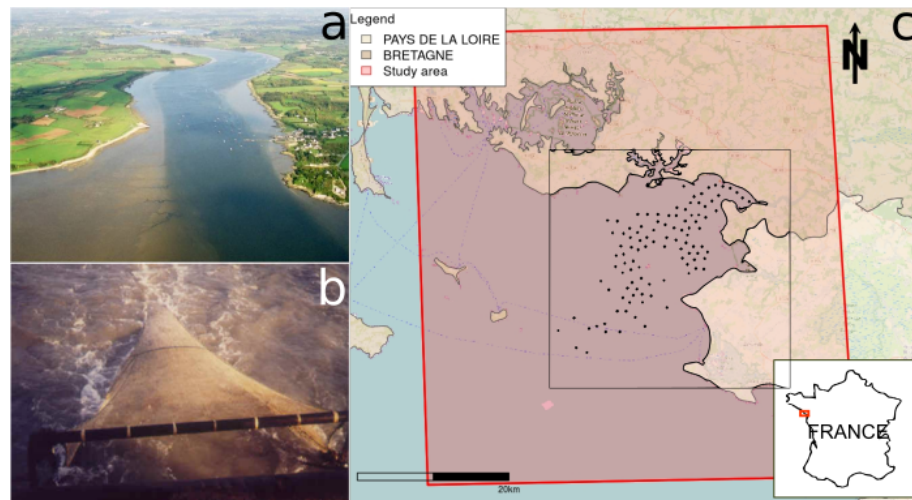
Déterminer les facteurs ayant une influence sur la distribution des poissons plats (*Solea solea*) en Baie de Vilaine et cartographier la distribution moyenne des densités.

- 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

## Données

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

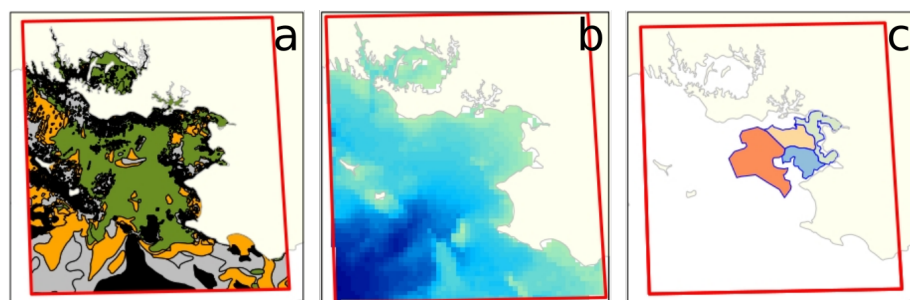
- 1984 – 2010
- En automne
- Juvéniles de l'année (Âge 0)
  - Nb individus / 1000m<sup>2</sup>



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

## 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.

## 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

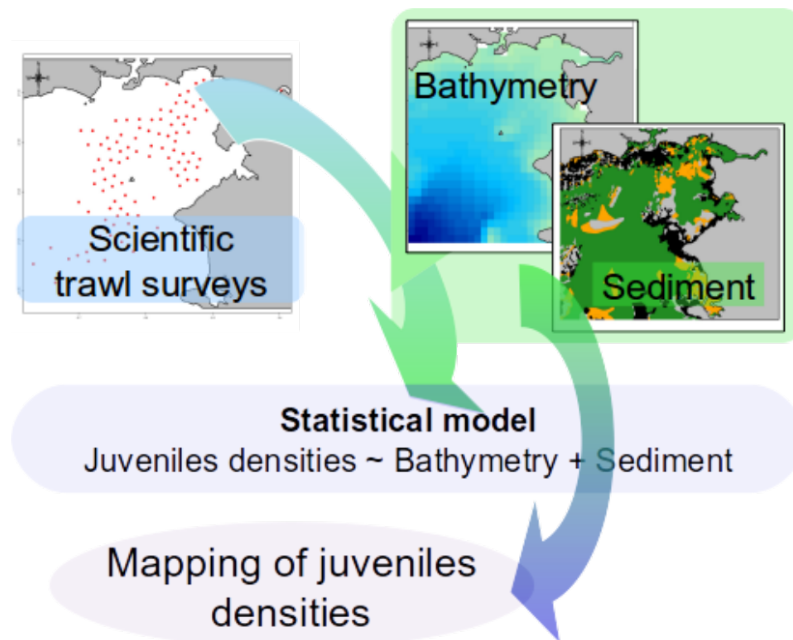


Figure 5 – Procédure pour un modèle de distribution d'espèce

## 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 ?*

## Préparation

### 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



## Débutons avec R

- Créer un projet Rstudio dans le dossier principal de travail.
- Ouvrez le script R : “Quick\_AllDataModel\_Teacher.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")
```

## Exploration des données

### Étapes

*Souvenez-vous : Définissez ce que vous cherchez, à quelles questions vous souhaiteriez répondre !*

- Explorer la répartition du plan d'échantillonnage en fonction des covariables environnementales
- Explorer les données d'observation au regard des covariables environnementales pour détecter de potentielles corrélations
- Explorer les interactions entre les effets des covariables sur les observations
- Explorer les lois de distribution possibles (gaussien, log-normal, ...) des observations en fonctions des combinaisons de covariables

Les scripts qui sont fournis ne sont que des exemples, ils ne sont pas des solutions ! Faites vos propres tests !

### Liste des différentes étapes

- Lire le jeu de données spatialisé (Fig. 6)
- Ajouter une nouvelle covariable : la bathymétrie divisée en classes
  - "< 5 m", "5-10 m", "10-20 m", "20-50 m"
- Explorer la répartition des observations en fonctions des covariables
  - Centrer l'analyse sur l'année, la bathymétrie et le sédiment
  - *Que remarquez-vous ?*
- Explorer les covariables ayant potentiellement des effets sur les densités
  - *Quelles covariables pourraient avoir une influence ?* (Fig. 7)

Les modèles statistiques que nous allons utiliser peuvent se résumer de cette façon :

$$Density = Covar1 + Covar2 + Noise$$

Comme vous le savez, on cherche toujours à savoir si les données sont gaussiennes pour pouvoir procéder à l'analyse statistique. Si elles ne sont pas gaussiennes, nous devons définir le type de distribution pour pouvoir utiliser une transformation de données.

- Explorer la distribution des données
- *Quelle est la distribution la plus intéressante ?*

La figure 8 montre différents exemples de distributions.

Exemple de l'effet de deux facteurs (Fig. 9)

- *Qu'en pensez-vous ?*
- Explorer les interactions potentielles entre les covariables

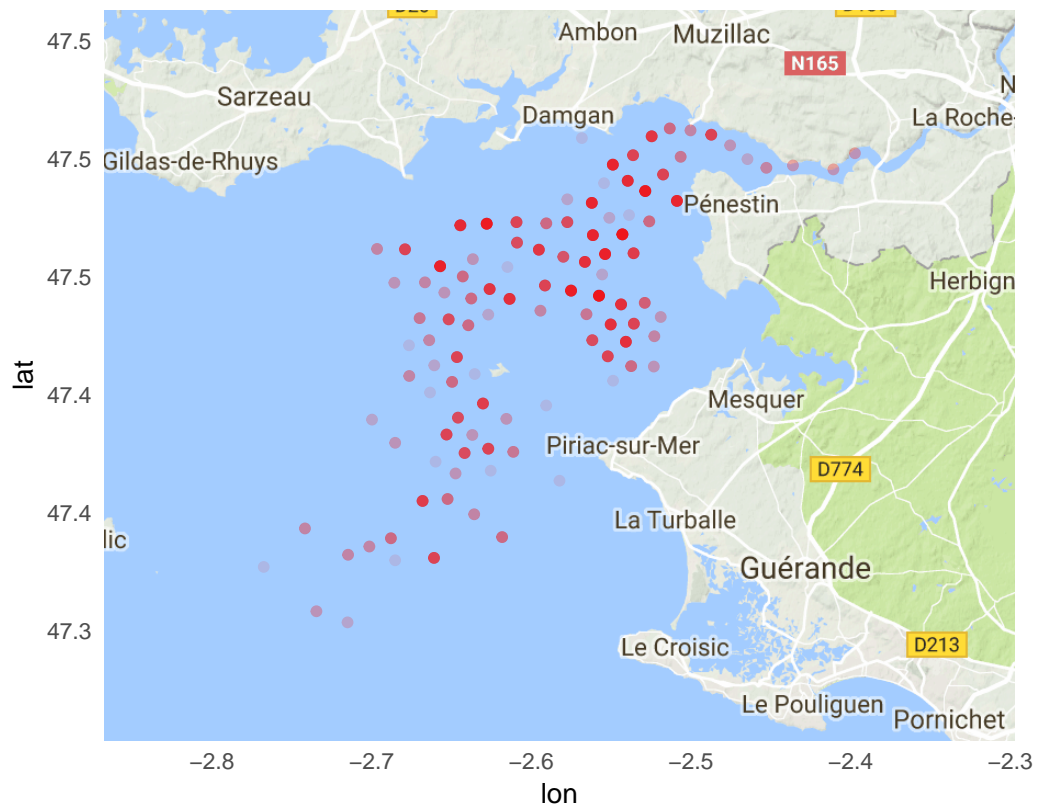


Figure 6 – Répartition des stations d'échantillonnage

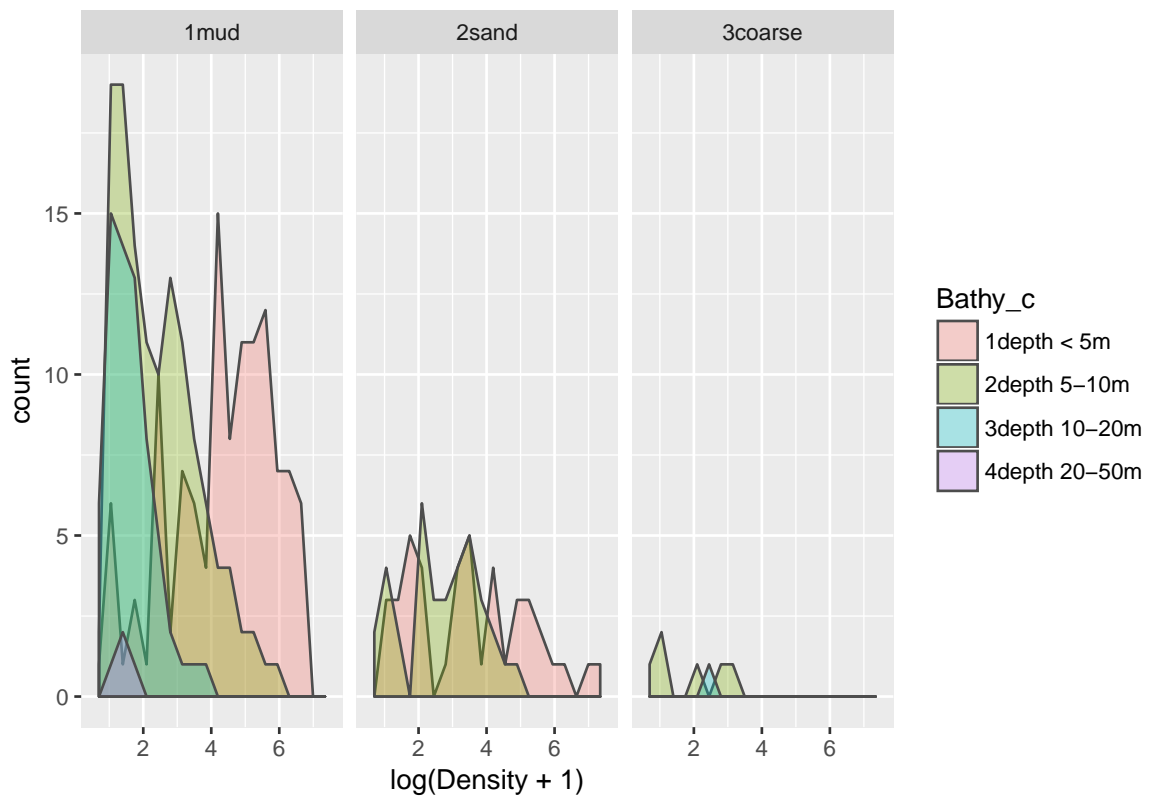


Figure 7 – Densités (log-transformées) en fonction de la bathymétrie et des sédiments

- *Qu'en pensez-vous ?*
- Comme aide à l'interprétation, utiliser l'exemple théorique de la figure 10

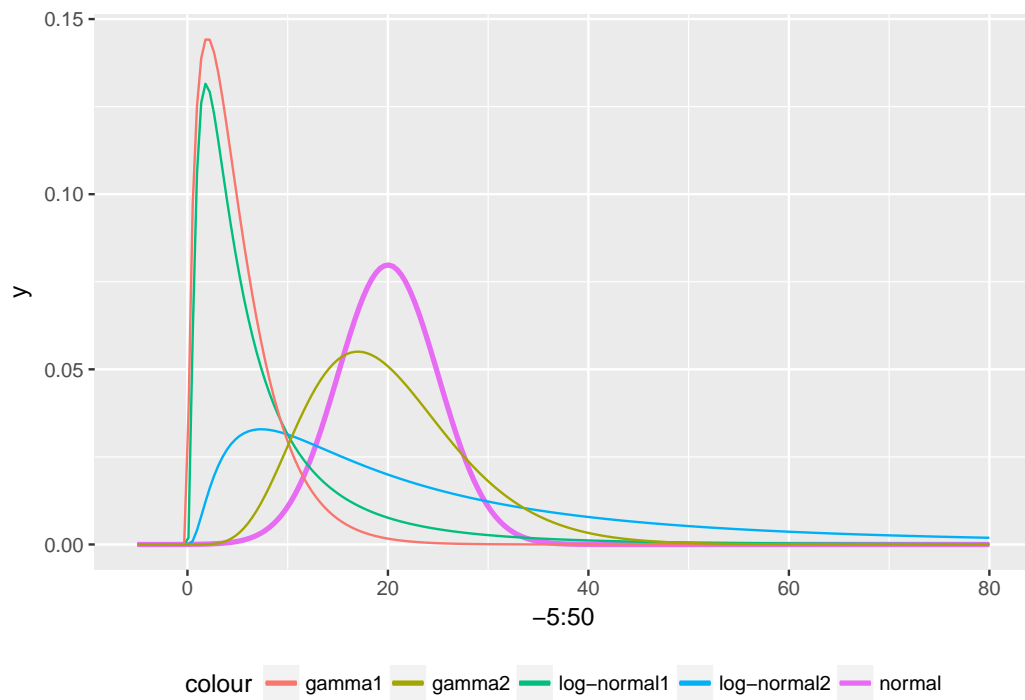


Figure 8 – Différents exemples de distributions

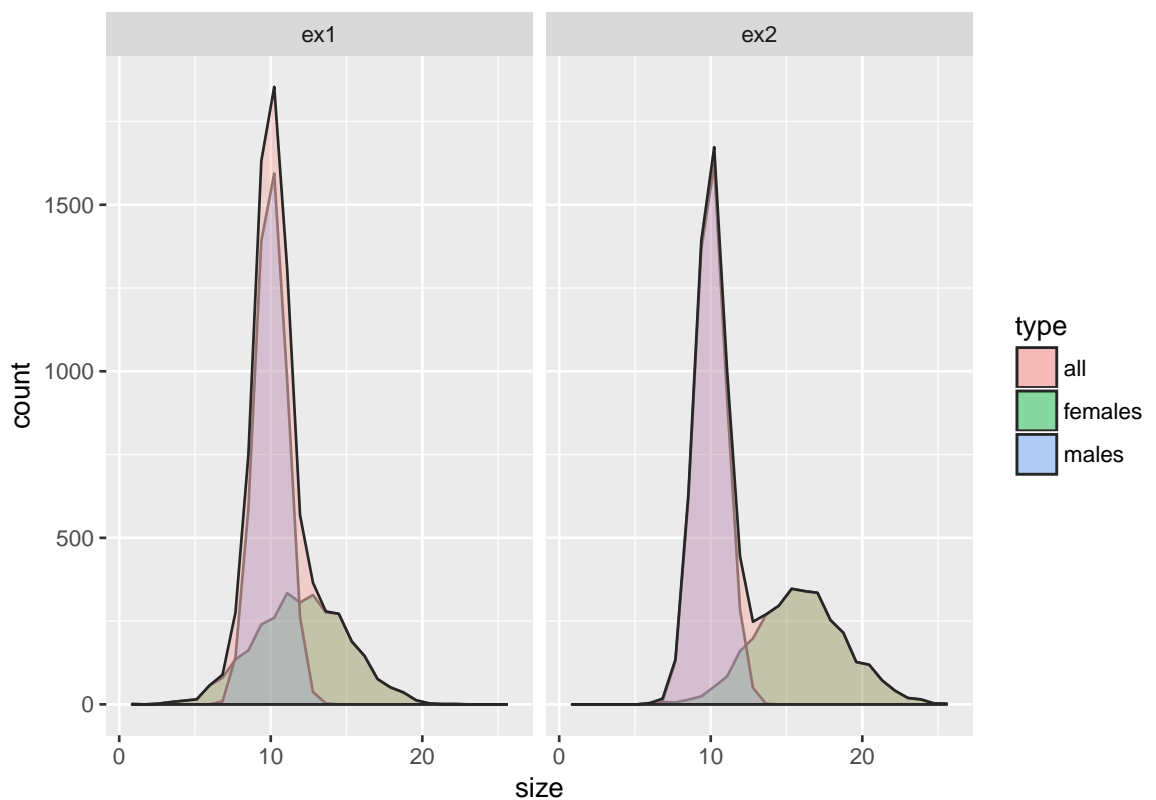


Figure 9 – Différents effets de deux facteurs

## Modélisation

### Étapes

*Souvenez-vous: Définissez ce que vous cherchez, à quelles questions vous souhaiteriez répondre !*

- Tester les différentes formes de modèles au regard des combinaisons de covariables et des



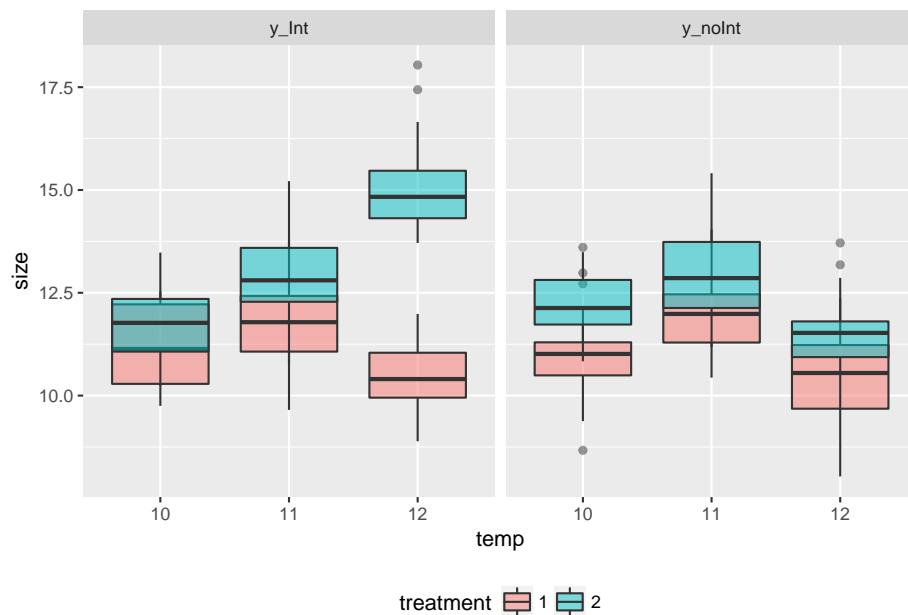


Figure 10 – Interaction entre traitements et temp  rature

formes de distributions des r  sidus

- Comparer les mod  les    l’aide des outils statistiques    disposition (AIC, anova, ...), de la validation crois  e mais aussi de la connaissance du jeu de donn  es et des questions cibl  es
- Analyser les r  sidus des mod  les. Analyser leur distribution et s’assurer que les hypoth  ses de construction sont v  rifi  es.

*C’est seulement lorsque les hypoth  ses sur la distribution des r  sidus sont v  rifi  es, que les covariables et les interactions s  lectionn  es peuvent commencer      tre interpr  t  es...*

Les scripts qui sont fournis ne sont que des exemples, ils ne sont en aucun cas les meilleures solutions ! Fa  tes vos propres tests !

## Interpr  ter les sorties de mod  les

Lorsque vous ajustez un mod  le lin  aire (lm ou glm), vous pouvez utiliser diff  rents tests statistiques et visuels qui r  pondent    diff  rentes questions. Votre question principale pourrait   tre :

- “Est-ce que mes covariables ont un effet sur mes observations ?”. En r  alit  , ce n’est pas exactement la question    laquelle va r  pondre votre mod  le. Ce serait plut  t “Est-ce que les covariables que j’ai utilis  es expliquent une part de la variabilit   de mes observations ?”

Pour que vous puissiez interpr  ter les diff  rentes sorties de mod  les, dans ce document, nous allons regarder le mod  le suivant :

`lm(Density Bathy + Sedim, data = dataset)`

*Ce mod  le n’est pas forc  ment le meilleur mod  le    choisir !*

### Summary(lm)

Cette fonction montre un tableau de tests de significativit   (Table 1). Ce sont des tests de Student. Ils testent si la valeur estim  e pour un effet est ou non significativement diff  rente de z  ro. Ainsi, si une covariable a un effet non significativement diff  rent de l’effet nul, il est probablement inutile de la conserver dans le mod  le.

- Ici, ce qui est appel   “Intercept” est l’effet de base. Dans une   quation  $y = a.x + b$ , l’“intercept” serait  $b$ . Ici, c’est un peu diff  rent car il y a des covariables au format facteur (“Sedim”). Dans cet exemple, l’“intercept” montre une “p-value” proche de z  ro, ce qui signifie que son effet (`estimate ~ 100`) est significativement diff  rent de z  ro.

Table 1 – Exemple d’une sortie de ‘summary(lm)’

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	173.4	14.27	12.152	0.000
Bathy	14.5	1.71	8.472	0.000
Sedim2sand	-19.5	19.07	-1.020	0.308
Sedim3coarse	-39.4	57.21	-0.689	0.491

- La covariable “Bathy” est continue. Dans une   quation de type  $y = a.x + b$  ce serait  $a$ . Dans cet exemple, son effet `estimate`  $\sim 6$ , est significativement diff  rent de z  ro (`p-value`  $\sim 0$ ).
- La covariable “Sedim” est un facteur    trois niveaux. Dans le “summary”, vous ne pouvez en voir que deux (“2sand”, “3coarse”). En r  alit  , les effets des niveaux de facteur (`estimate`) sont compar  s au premier niveau (“1mud”), pour lequel l’effet est   gal    z  ro. Dans l’  quation  $Y = a * Bathy + b[Sedim] + c$ , l’estimation de la moyenne  $Y$  pour `Sedim = "1mud"` est  $Y = a*Bathy + 0 + c$ . Dans les r  sultats du tableau, `b["2sand"]` est donc non significativement diff  rent de `b["1mud"] = 0`, et, `b["3coarse"]` avec une `p-value` = 0.07, n’est pas non plus significativement diff  rent de `b["1mud"] = 0`.

### Analyse de r  sidus

Une hypoth  se de construction d’un mod  le lin  aire est l’homosc  dasticit  , aussi appel  e homog  n  it   de la variance. Cela signifie que la variance de la r  ponse  $y$  est la m  me quelque soit la valeur du pr  dicteur  $x$ . Dans un mod  le Gaussien classique ajustant  $x$      $y$  ainsi:  $y = a.x + b + \epsilon$ , la variable  $\epsilon$  repr  sente les r  sidus du mod  le. Ils sont suppos  s   tre centr  s sur z  ro et avec une variance Gaussienne, leur distribution suivant ainsi la loi Gaussienne  $\epsilon \sim N(0, \sigma)$

Lorsqu’on simule un tel mod  le, par exemple  $y = 2.x + 5 + \epsilon$ , on observe la figure 11, avec une homog  n  it   de la distribution des observations autour de l’ajustement, et une distribution Gaussienne des r  sidus comme le montrent l’histogramme et le “qqplot”.

```
# Example of homogeneous residuals
n <- 1000
epsilon <- rnorm(n, 0, 5)
x <- runif(n, 0, 10)
y <- 2*x + 5 + epsilon

par(mfrow = c(1,3))
plot(x, y, pch = 20)
abline(5, 2, col = "red", lwd = 2)
hist(epsilon, breaks = 20, col = "grey")
qqnorm(epsilon); qqline(epsilon, col = "red")
```

   partir de cet exemple, vous pouvez d  finir les diagnostics graphiques n  cessaires pour v  rifier vos hypoth  ses de construction de mod  le. Lorsqu’on utilise le m  me mod  le que pr  c  demment (Fig.12) :

- **Residuals vs Fitted** - Dans cet exemple, on peut voir que la variabilit   des r  sidus augmente avec les valeurs pr  dites, ce qui va    l’encontre de l’homog  n  it   de la variance.
- **Scale-Location** est en accord avec la figure pr  c  dente car on voit une augmentation de la deviance des r  sidus quand les pr  dictions augmentent.
- **Normal Q-Q** - Cette figure appel  e “qqplot” montre la divergence entre les quantiles th  oriques d’une loi Gaussienne et les quantiles r  els de la distribution des r  sidus du mod  le. La divergence est importante pour les valeurs   lev  es, ce qui montre une queue de distribution plus longue qu’une loi Normale.
- **Hist. of residuals** - L’histogramme des r  sidus est en accord avec le qqplot car on voit clairement une distribution qui n’est pas une Gaussienne centr  e sur z  ro. Cette distribution a une longue queue de distribution avec beaucoup plus de valeurs positives que de valeurs n  gatives.

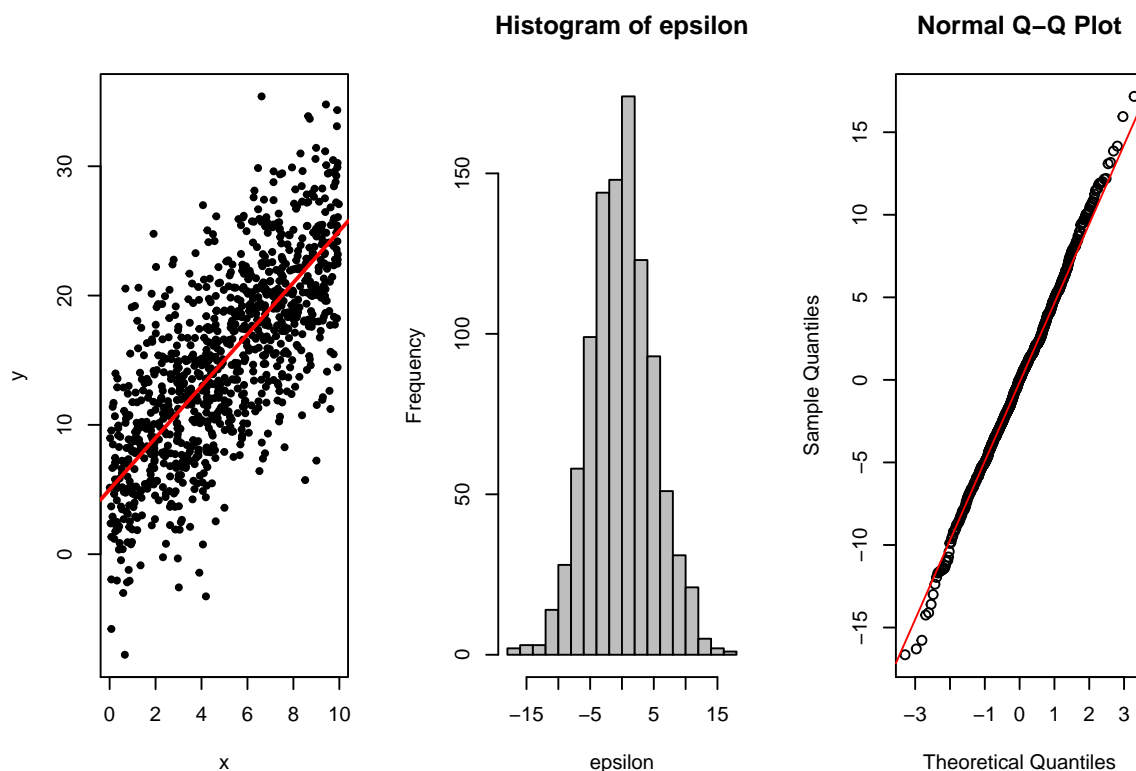


Figure 11 – Simulation d'une relation linéaire entre  $x$  et  $y$  avec un résidu Gaussien.

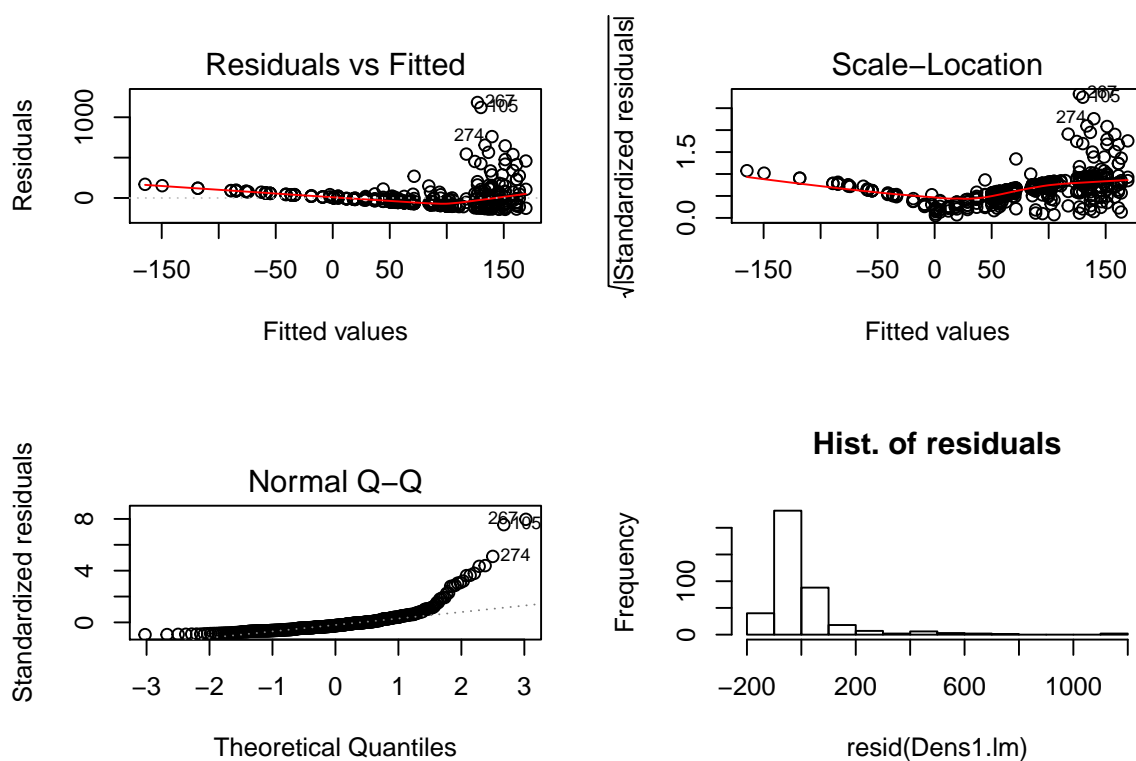


Figure 12 – Figures de diagnostic d'un modèle linéaire permettant de vérifier les hypothèses de construction.

## Analyse de variance

La question à laquelle répond une **anova** (avec un test du Chi-2) est : Est-ce que la covariable ajoutée augmente significativement la vraisemblance du modèle (ou a réduit la déviance résiduelle), comparé au modèle précédent, sans cette covariable ?

Table 2 – Exemple d’une ‘sortie’ de `anova(lm)`

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
Bathy	1	1607692	1607692	71.983	0.000
Sedim	2	32206	16103	0.721	0.487
Residuals	397	8866764	22334	NA	NA

Table 3 – Comparaison de la déviance expliquée par différents modèles avec un nombre croissant de paramètres

model	Res.Df	RSS	Df	Sum of Sq	Pr(>Chi)	"  "	df	AIC
lm1	8	3192	NA	NA	NA		3	92.0
lm2	6	1272	2	1920	0.003		5	86.8
lm3	4	645	2	626	0.144		7	84.1

- NULL est le test pour un modèle sans covariable, c’est le modèle qui estime la moyenne :  $Density \sim constant$ .
- Bathy est le test pour un modèle uniquement avec la Bathy :  $Density \sim Bathy$ . La p-value est proche de zéro, indiquant que le gain de déviance expliquée en ajoutant la Bathy est significativement différent de zéro.
- Sedim est le test pour un modèle avec Bathy et Sedim, dans cet ordre :  $Density \sim Bathy + Sedim$ . La p-value  $\sim 0.2$  indique qu’il y a un risque de 20% que la déviance expliquée en ajoutant la covariable Sedim au modèle contenant déjà la Bathy soit nulle.

### Critères d’Akaike (AIC) et Bayesian (BIC)

L’AIC et le BIC sont des critères de qualité d’ajustement pénalisés par le nombre de paramètres estimés. La description (traduite) de ces fonctions dans R est :

Function générique calculant “Le Critère d’Information” d’Akaike pour un ou plusieurs modèles ajustés pour lesquels une “log-vraisemblance” peut être obtenue, en utilisant la formule  $IC = -2 * \log - likelihood + k * npar$ , où `npar` représente le nombre de paramètres estimés, et `k` = 2 pour l’AIC classique, ou `k` =  $\log(n)$  (`n` étant le nombre d’observations) pour le BIC ou SBC (Schwarz’s Bayesian criterion).

En effet, plus vous ajoutez de paramètres dans un modèle, plus vous avez de chances que le modèle s’ajuste parfaitement aux données (Table 3, Fig. 13). L’AIC diminue avec la déviance résiduelle et augmente avec le nombre de paramètres ajustés. Plus l’AIC est bas, plus parsimonieux est le modèle.

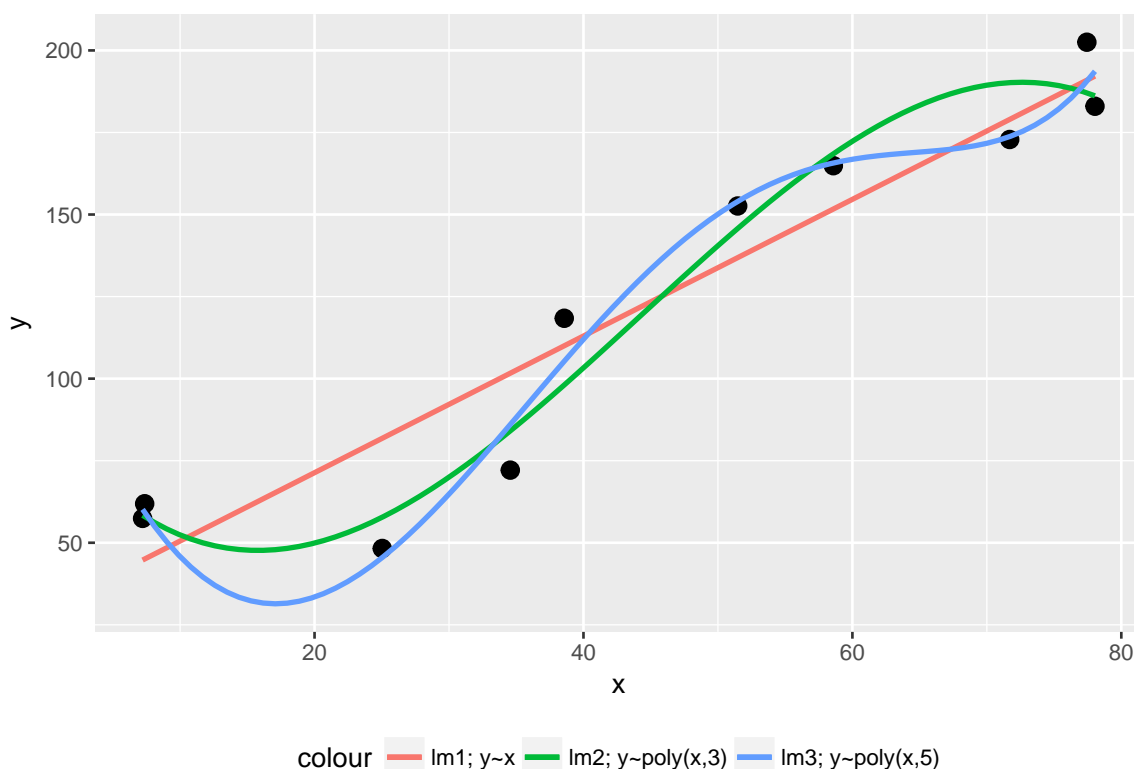
### Trouver le meilleur modèle

La fonction `lm` n’est utilisée que pour des modèles avec une distribution Gaussienne des résidus. Pour tester d’autres types de distributions, il faut utiliser `glm`, avec un paramètre pour la famille de distribution (`family`). Vous pouvez utiliser des distributions qui autorisent une plus grande queue de distribution que la loi Normale. Parmi les familles disponibles, vous pouvez tester `poisson`, `quasipoisson`, `Gamma`, `Log-gaussian`.

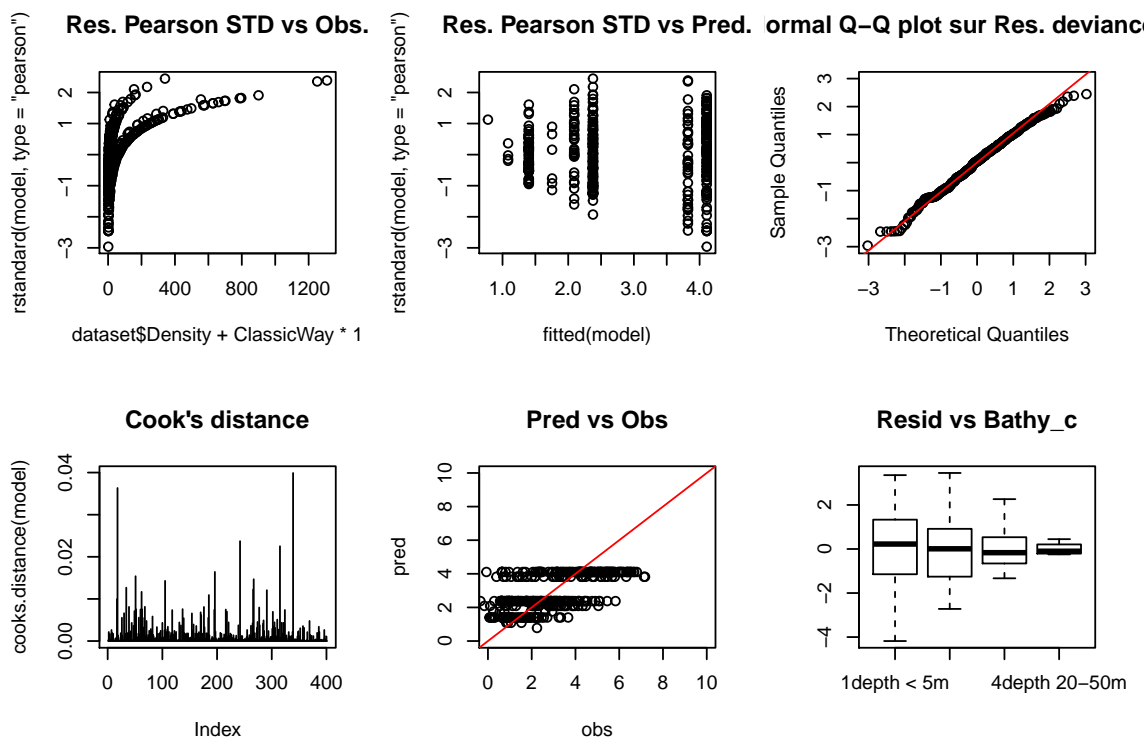
- *Quel est le meilleur modèle au regard des différents critères évoqués ?*

### Exploration des sorties du meilleur modèle

- *Que pouvez-vous dire sur le diagnostic complet de votre modèle ?* (Fig. 14)



**Figure 13** – Différents modèles ajustés sur les mêmes données mais avec un nombre de paramètres ajustés croissant.



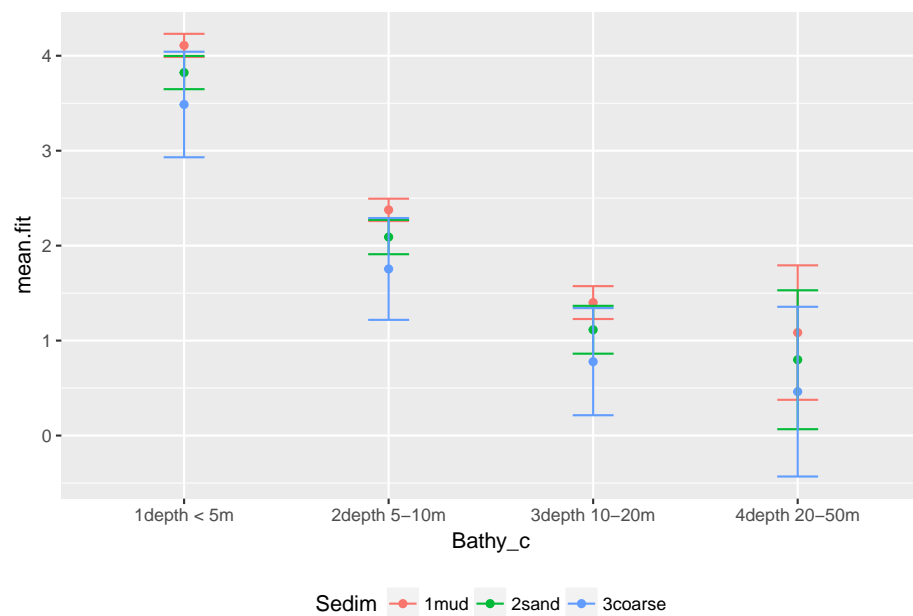
**Figure 14** – Diagnostic du meilleur GLM sélectionné

## Prédictions du modèle

Lorsque vous êtes satisfaits du modèle sélectionné, vous pouvez faire des prédictions

- Utiliser le fichier csv fourni pour voir l'effet des covariables sélectionnées (Fig. 15)





**Figure 15** – Pr  dictions du meilleur GLM s  lectionn  