



J.-M. Galharret,  
A. Philippe,  
N. Mercier

<https://galharret.github.io/WEBSITE>

# Modèle hiérarchique pour la détection de valeurs aberrantes

Applications à la datation en archéologie

19 Décembre 2023

# Plan

---

## 1 Introduction

Problématique

Event Model

Loi de shrinkage

## 2 Méthodologie générale

Modèle hiérarchique

$\theta_1, \dots, \theta_n$  échangeables

Outliers

Règle de décision

Stratégies de ré-estimation

## 3 Modèle hiérarchique

gaussien

Simulations

Exemple

Calibrage de  $\alpha$

## 4 Application 1 en archéologie

Tel Quasile

Comparaison

## 5 Autre modélisation : âge OSL

Relation fondamentale

Ajustement de  $d$

Modélisation

Estimation de  $s_0^2$

Estimation sur données réelles

Validation finale du modèle

Conclusion

En archéologie, quelle que soit la méthode de datation utilisée (C14,OSL,...), on est confronté au problème des outliers

- Erreur de mesure (laboratoire),
- Erreur de prélèvement (fouilles archéologiques).

Logiciels de modélisation chronologique

- OxCal : Modèle de Bronk Ramsey (2009),
- Chronomodel Lanos and Philippe (2017-18).

Notre approche :

- Identification des outliers via le modèle robuste,
- Ré-estimation du paramètre sur le sous-échantillon.

Modèle hiérarchique  $\rightsquigarrow$  estimer l'âge d'un évènement à partir de la datation de  $n$  objets le caractérisant.

- Pour le  $i$ -ème objet d'âge  $A_i$ , le laboratoire fournit la mesure  $X_i$  avec une erreur  $s_i$

$$X_i | A_i, s_i \sim \mathcal{N}(A_i, s_i^2).$$

- On suppose que les âges  $A_1, \dots, A_n$  sont contemporains de  $A$

$$A_i | A, \sigma \sim \mathcal{N}(A, \sigma^2).$$

.



- Event model<sup>1</sup>

$$A_i | A, \sigma_i \sim \mathcal{N}(A, \sigma_i^2), \quad \sigma_i^2 \text{ i.i.d.}$$

1. Lanos and Philippe (2017) Hierarchical Bayesian modeling for combining dates in archaeological context. Journal de la Société Française de Statistique, 158(2) :72(88)

Loi de shrinkage uniforme pour les  $\sigma_i^2$

$$\frac{s_0^2}{s_0^2 + \sigma_i^2} \stackrel{i.i.d}{\sim} \mathcal{U}[0, 1],$$

où  $s_0^{-2} = \sum_{i=1}^n s_i^{-2}$

Avantages de ce choix :

- La médiane de  $\sigma_i^2$  est égale à  $s_0^2$ .  $\rightsquigarrow$  poids identique entre les erreurs de mesure et de modèle.
- Cette loi a des moments infinis ( $\rightsquigarrow$  outliers).

Inconvénient : Galharret & al. (2021) ont montré que la loi a posteriori a également des moments infinis.

# Plan

---

- 1** Introduction
  - Problématique
  - Event Model
  - Loi de shrinkage
- 2** Méthodologie générale
  - Modèle hiérarchique
  - $\theta_1, \dots, \theta_n$  échangeables
  - Outliers
  - Règle de décision
  - Stratégies de ré-estimation
- 3** Modèle hiérarchique gaussien
  - Simulations
  - Exemple
  - Calibrage de  $\alpha$
- 4** Application 1 en archéologie
  - Tel Quasile
  - Comparaison
- 5** Autre modélisation : âge OSL
  - Relation fondamentale
  - Ajustement de  $d$
  - Modélisation
  - Estimation de  $s_0^2$
  - Estimation sur données réelles
  - Validation finale du modèle
  - Conclusion

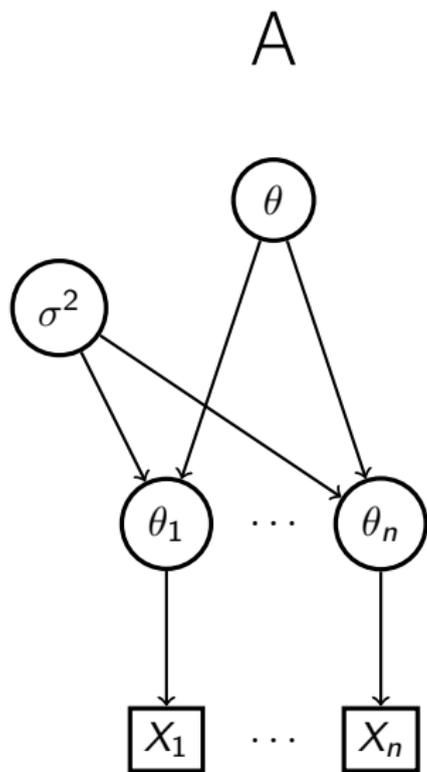
$X_1, \dots, X_n$   $n$  observations de vraisemblance

$$X_1, \dots, X_n \sim p^{(n)}(\cdot \mid \theta) \quad (1)$$

Un modèle hiérarchique classique (Spiegelhalter (2004), Congdon (2002), Gelman (2008))

$$\begin{aligned} p^{(n)}(X_1, \dots, X_n \mid \theta) &= \int f^{(n)}(X_1, \dots, X_n, \theta_1, \dots, \theta_n \mid \theta) d\theta_1 \dots d\theta_n \\ &= \int f^{(n)}(X_1, \dots, X_n \mid \theta_1, \dots, \theta_n, \theta) \\ &\quad \times \pi_1(\theta_1, \dots, \theta_n \mid \theta) d\theta_1 \dots d\theta_n. \end{aligned}$$

$\rightsquigarrow$  Les variables aléatoires  $|\theta_1 - \theta|, \dots, |\theta_n - \theta|$  mesurent l'hétérogénéité entre  $X_1, \dots, X_n$ .

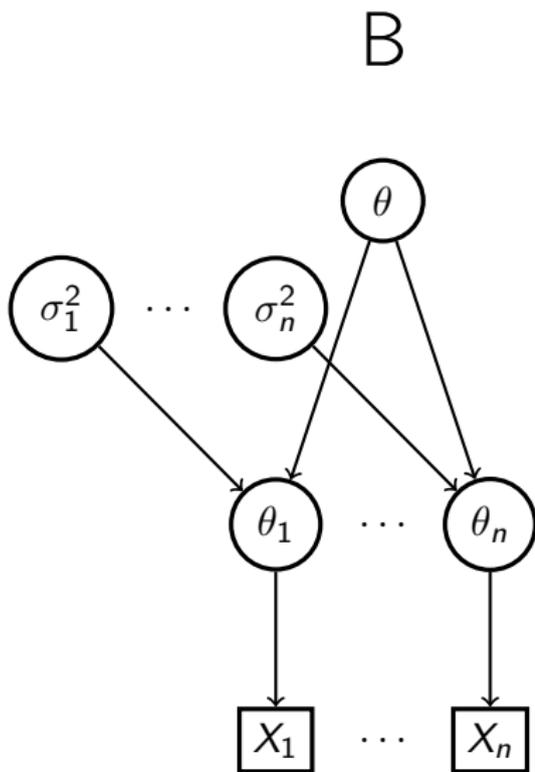


$$\pi_1(\theta_1, \dots, \theta_n | \theta) = \prod_{i=1}^n \pi_{\sigma^2}(\theta_i | \theta),$$

$$f^{(n)}(X_1, \dots, X_n | \theta_1, \dots, \theta_n, \theta) = \prod_{i=1}^n f(X_i | \theta_i),$$

$$\sigma^2 = \text{Var}(\theta_i - \theta | \theta)$$

$$= \text{Var}(\theta_i | \theta).$$



$$\sigma_i^2 = \text{Var}(\theta_i | \theta),$$
$$\pi_1(\theta_1, \dots, \theta_n | \theta) = \prod_{i=1}^n \pi_{\sigma_i^2}(\theta_i | \theta).$$

Basée<sup>2</sup> sur la comparaison des lois a priori et a posteriori des  $\sigma_i^2$   
Soit  $\alpha \in ]0, 1[$  et soit  $q_{1-\alpha}$  le quantile d'ordre  $1 - \alpha$  de la loi a priori de  $\sigma_i^2$

$$\mathbb{P}(\sigma_i > q_{1-\alpha}) = \alpha,$$

l'observation  $X_j$  est un outlier si

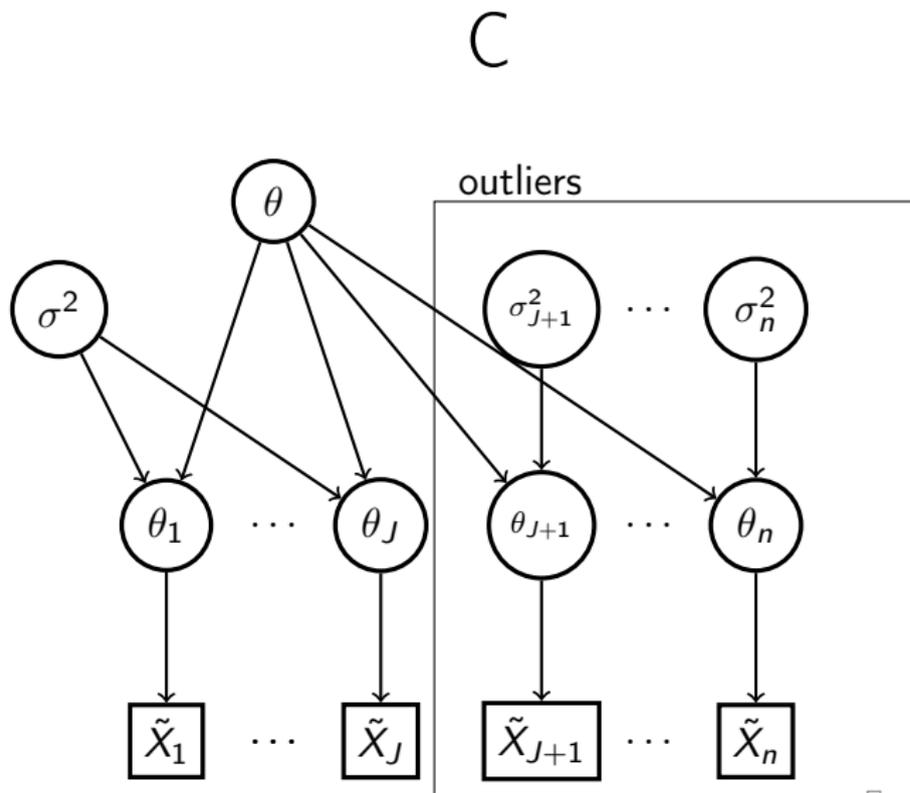
$$\mathbb{P}(\sigma_i > q_{1-\alpha} | X_1, \dots, X_n) \geq \alpha,$$

Dans la suite on note  $(\tilde{X}_i)_{i \in \{1, \dots, J\}}$  le sous-échantillon de  $(X_i)_{i \in \{1, \dots, n\}}$   
d'où ont été exclus les outliers détectés  $(\tilde{X}_i)_{i \in \{J+1, \dots, n\}}$ .

---

2. Galharret, J-M, Philippe, A., & Mercier, N. (2021). Detection of outliers with a Bayesian hierarchical model : application to the single-grain luminescence dating method. Electronic Journal Of Applied Statistical Analysis.

## Stratégies de ré-estimation



# Plan

---

- 1** Introduction
  - Problématique
  - Event Model
  - Loi de shrinkage
- 2** Méthodologie générale
  - Modèle hiérarchique
  - $\theta_1, \dots, \theta_n$  échangeables
  - Outliers
  - Règle de décision
  - Stratégies de ré-estimation
- 3** **Modèle hiérarchique gaussien**
  - Simulations
  - Exemple
  - Calibrage de  $\alpha$
- 4** Application 1 en archéologie
  - Tel Quasile
  - Comparaison
- 5** Autre modélisation : âge OSL
  - Relation fondamentale
  - Ajustement de  $d$
  - Modélisation
  - Estimation de  $s_0^2$
  - Estimation sur données réelles
  - Validation finale du modèle
  - Conclusion

Soit  $\tau \in \{5\%, 10\%, 20\%\}$  le taux de contamination et  $\mu \in \{5, 10, 20\}$  le décentrage. On simule  $n$  observations  $X_1, \dots, X_n$  dont

- $J = \lceil n\tau \rceil$  vraies observations  $X_1, \dots, X_J$

$$X_i \sim \mathcal{N}(\theta, s_i^2),$$

- $n - J$  outliers  $X_{J+1}, \dots, X_n$

$$X_i \sim \mathcal{N}(\theta + \mu, s_i^2).$$

On effectue  $B = 1000$  répliques.

# Modèle hiérarchique gaussien

## Exemple

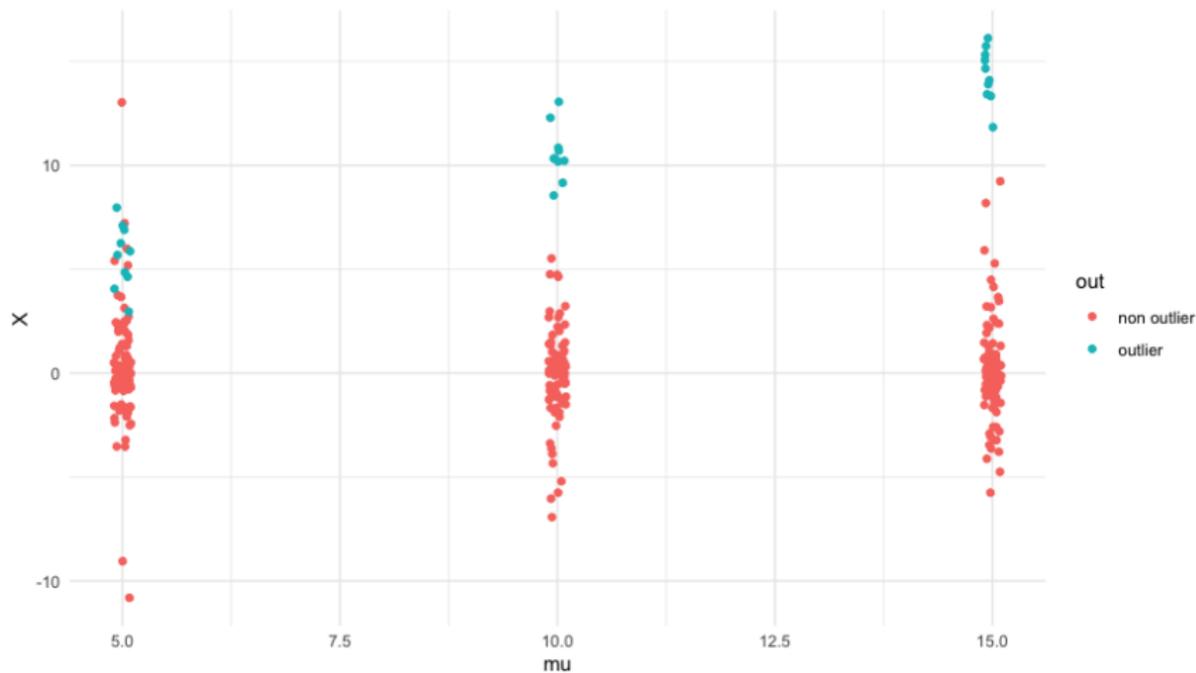
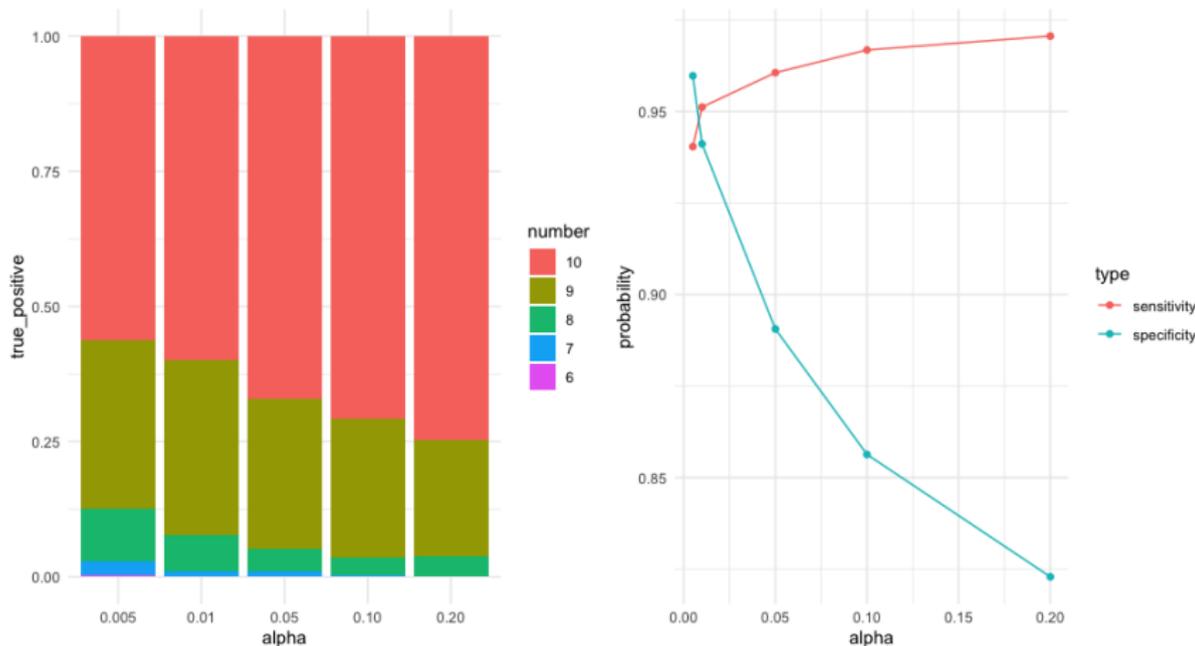


Figure – Exemple d'une simulation avec  $\tau = 0.10$

# Modèle hiérarchique gaussien

## Calibrage de $\alpha$



**Figure** – Variation du nombre de vrais positifs (gauche) et de la sensibilité/spécificité (droite) en fonction de  $\alpha$  pour  $\tau = 10\%$  et  $\mu = 10$ .

# Modèle hiérarchique gaussien

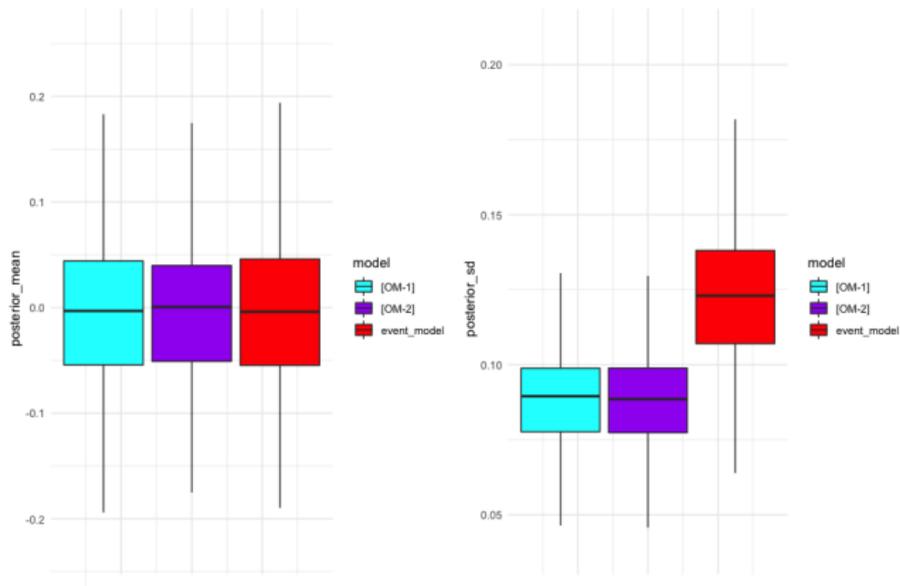
## Résultats des simulations

	$\tau \backslash \mu$	15	10	5
Sensitivity	5%	0.99	0.96	0.81
	10%	0.99	0.96	0.80
	20%	0.99	0.97	0.80
Specificity	5%	0.89	0.89	0.89
	10%	0.89	0.89	0.89
	20%	0.89	0.89	0.89

**Table** – Estimation of the sensitivity and the specificity as function of contamination rate  $\tau$  and the mean value  $\mu$  of the distribution of the outliers. The cut-off is fixed to  $\alpha = 0.05$ , and the number of replications is  $N = 1000$ .

# Modèle hiérarchique gaussien

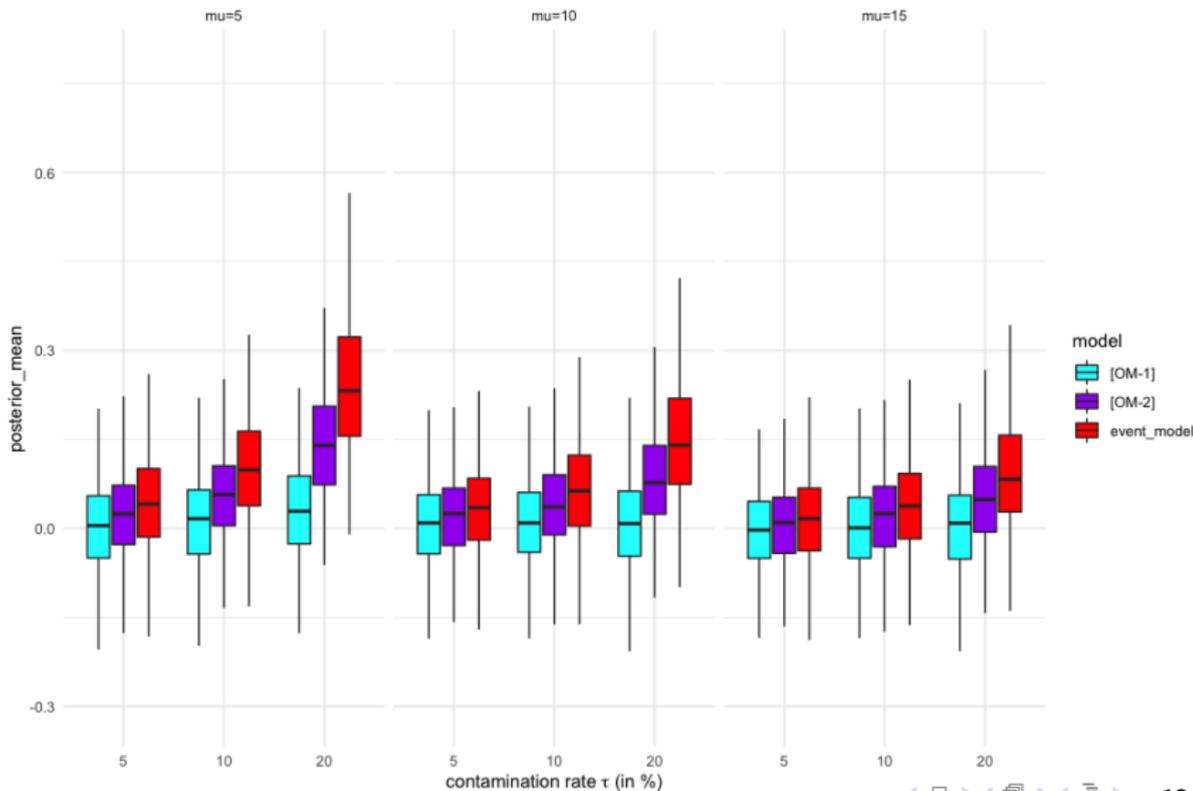
## Comparaison sous $H_0$

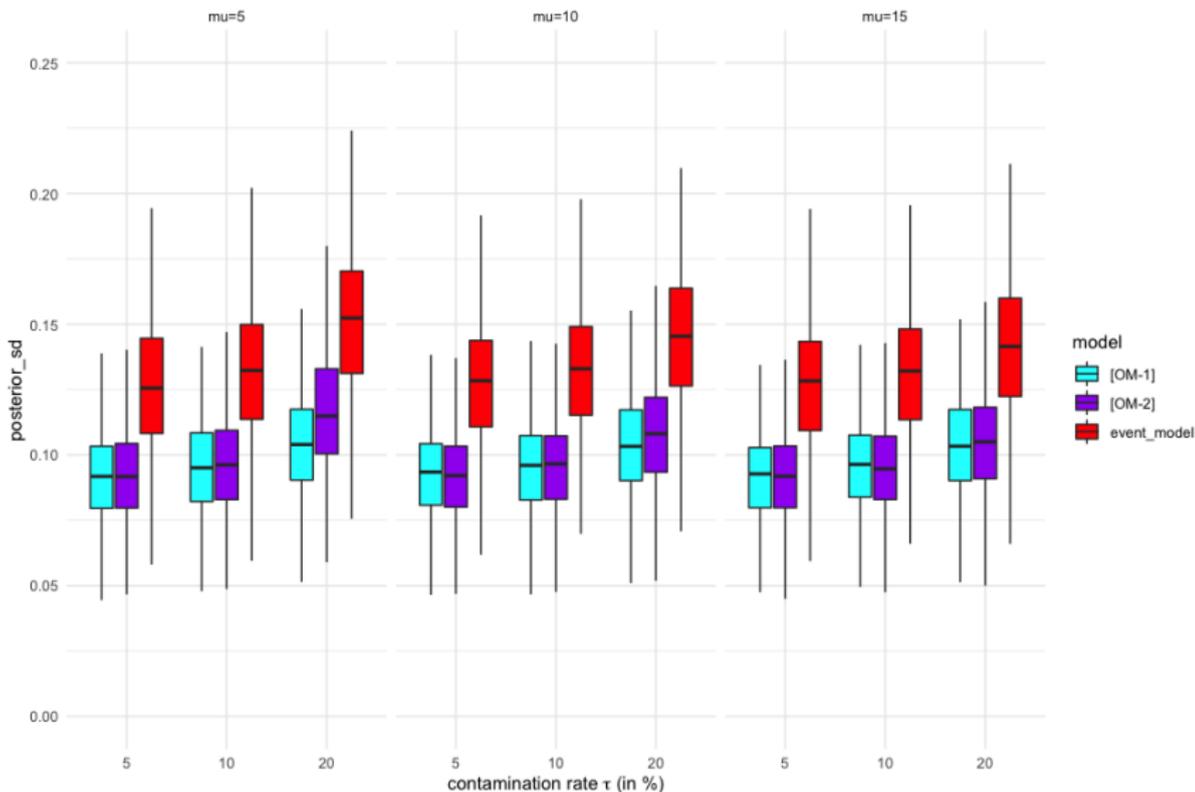


**Figure** – Comparison of the three following models : [OM-1], [OM-2] and the event model on simulated dataset without outlier ( $\tau = 0$ ). We represent the boxplot of the mean (left) and standard deviation (Right) of the posterior distribution of  $\theta$ .

# Modèle hiérarchique gaussien

## Comparaison sous $H_1$





**Figure** – Comparison of the three models [OM-1], [OM-2], event model on simulated dataset with outliers for different values of contaminated rates  $\tau$  and parameter  $\mu$ .

# Plan

---

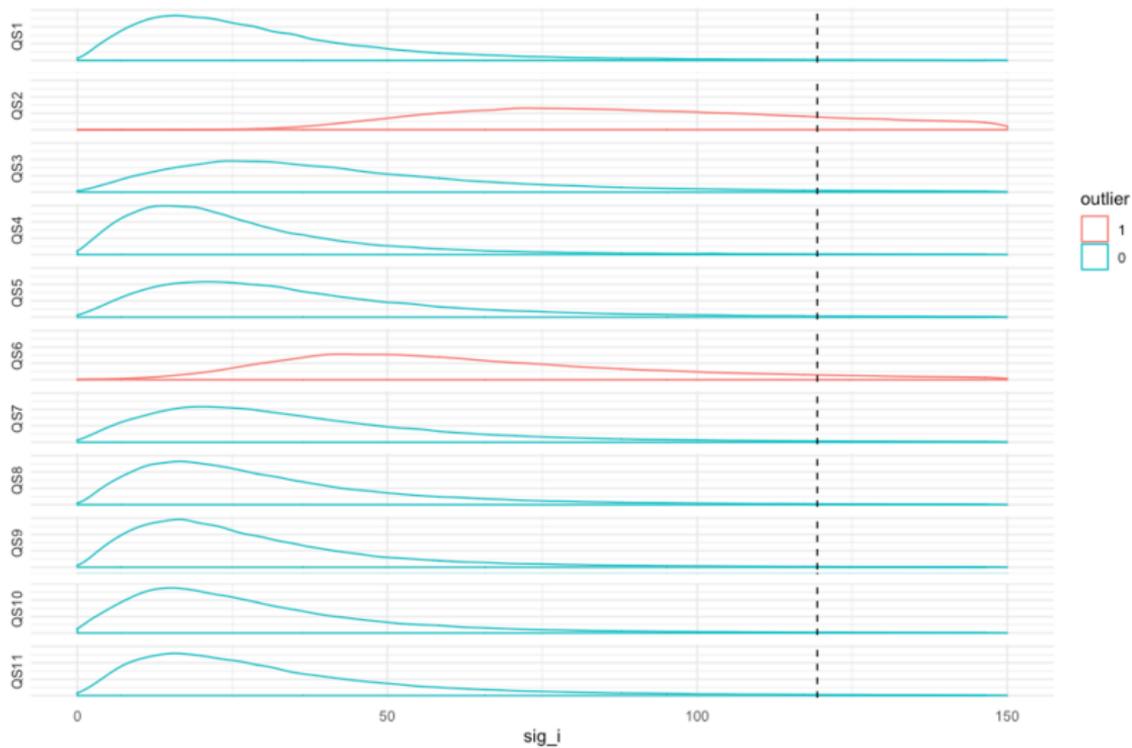
- 1** Introduction
  - Problématique
  - Event Model
  - Loi de shrinkage
- 2** Méthodologie générale
  - Modèle hiérarchique
  - $\theta_1, \dots, \theta_n$  échangeables
  - Outliers
  - Règle de décision
  - Stratégies de ré-estimation
- 3** Modèle hiérarchique gaussien
  - Simulations
  - Exemple
  - Calibrage de  $\alpha$
- 4** Application 1 en archéologie
  - Tel Quasile
  - Comparaison
- 5** Autre modélisation : âge OSL
  - Relation fondamentale
  - Ajustement de  $d$
  - Modélisation
  - Estimation de  $s_0^2$
  - Estimation sur données réelles
  - Validation finale du modèle
  - Conclusion

## Application 1 en archéologie

# Tel Quasile

- Tel Quasile est un site archéologique situé en Israël.
- Célèbre pour avoir contribué à la connaissance des Phillistins (installés entre le XII et X siècle avant JC)
- Bronk Ramsey a identifié les outliers dans les 12 échantillons datés par C14 en utilisant un mélange.





## Application 1 en archéologie

# Comparaison

method	ident	date $X_i$	$s_i$	$\mathbb{P}(\sigma_i > q_{.95} \mid X_1, \dots, X_n)$
$^{14}\text{C}$	QS1	2818	26	0.017
$^{14}\text{C}$	QS2	2692	24	<b>0.358</b>
$^{14}\text{C}$	QS3	2911	26	<b>0.046</b>
$^{14}\text{C}$	QS4	2853	25	0.016
$^{14}\text{C}$	QS5	2895	25	0.030
$^{14}\text{C}$	QS6	2753	22	<b>0.128</b>
$^{14}\text{C}$	QS7	2800	25	0.030
$^{14}\text{C}$	QS8	2882	28	0.020
$^{14}\text{C}$	QS9	2864	40	0.015
$^{14}\text{C}$	QS10	2818	38	0.019
$^{14}\text{C}$	QS11	2897	44	0.023

**Table** – Dates from Tell Qasile X and outputs of the decision rule : bold values indicate radiocarbon dates detected as outliers.

# Plan

---

## 1 Introduction

Problématique

Event Model

Loi de shrinkage

## 2 Méthodologie générale

Modèle hiérarchique

$\theta_1, \dots, \theta_n$  échangeables

Outliers

Règle de décision

Stratégies de ré-estimation

## 3 Modèle hiérarchique

gaussien

Simulations

Exemple

Calibrage de  $\alpha$

## 4 Application 1 en archéologie

Tel Quasile

Comparaison

## 5 Autre modélisation : âge OSL

Relation fondamentale

Ajustement de  $d$

Modélisation

Estimation de  $s_0^2$

Estimation sur données réelles

Validation finale du modèle

Conclusion

- La thermoluminescence permet de dater des matériaux archéologiques qui ont autrefois été chauffés.
- On date la dernière chauffe ou la dernière exposition à la lumière d'un échantillon contenant des minéraux de quartz.
- Méthode : [luminescence stimulée optiquement] l'énergie cumulée est alors libérée sous forme lumineuse et la quantité de lumière émise est proportionnelle au temps écoulé depuis la dernière chauffe.

↪ Outliers proviennent de grains mal blanchis (c'est à dire qui n'ont pas rejeté toute l'énergie cumulée).

## Autre modélisation : âge OSL

# Relation fondamentale

Dans le cadre de la datation de quartz par luminescence on estime son âge en utilisant la relation suivante

$$D \stackrel{\mathcal{L}}{=} A\dot{d}$$

où  $\dot{d}$  est le débit de dose,  $D$  est la dose équivalente absorbée et  $A$  l'âge cible.

- $\dot{d}$  n'est pas observable mais on sait simuler des échantillons de  $\dot{d}$  à une erreur systématique *int* prêt

$$\dot{d} = \tilde{\dot{d}}(1 + \text{int}\dot{\varepsilon}), \quad \dot{\varepsilon} \sim \mathcal{N}(0, 1).$$

- Première étape : on ajuste  $\dot{d}$  par un mélange gaussien de paramètres  $(K, p_1, \dots, p_K, \dot{\mu}_1, \dots, \dot{\mu}_K, \dot{\sigma}_1, \dots, \dot{\sigma}_K)$ .
- On a alors

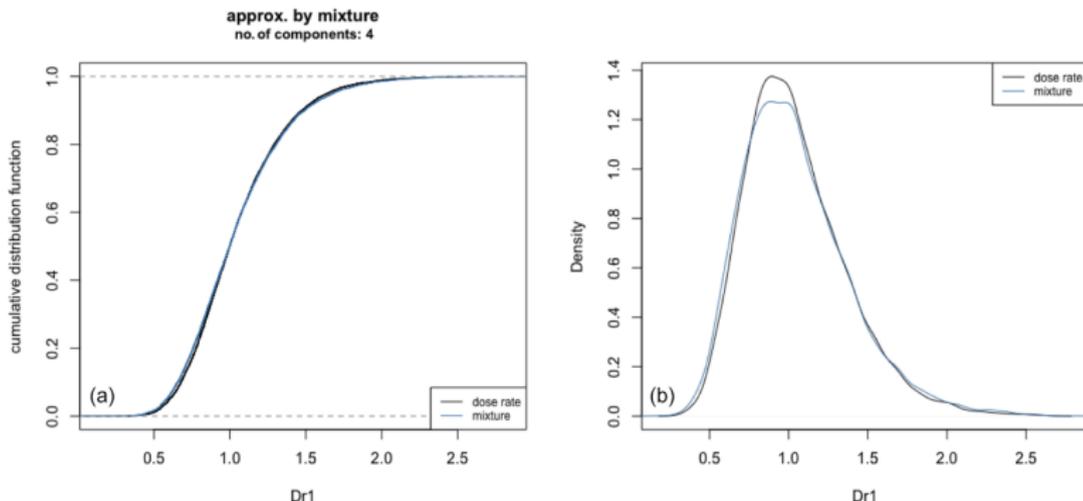
$$D \simeq \sum_{k=1}^K \dot{p}_k \mathcal{N}(A\dot{\mu}_k, A^2\dot{\sigma}_k^2).$$

## Autre modélisation : âge OSL

# Ajustement de $d$

- Pour  $K$  fixé, les paramètres du mélange sont estimés par l'algorithme EM.
- Pour le choix de  $K$ , on utilise  $BIC$ .

Figure – Exemple d'ajustement



$$\tilde{D}_j \sim \mathcal{N}(D_j, s_{D_j}^2)$$

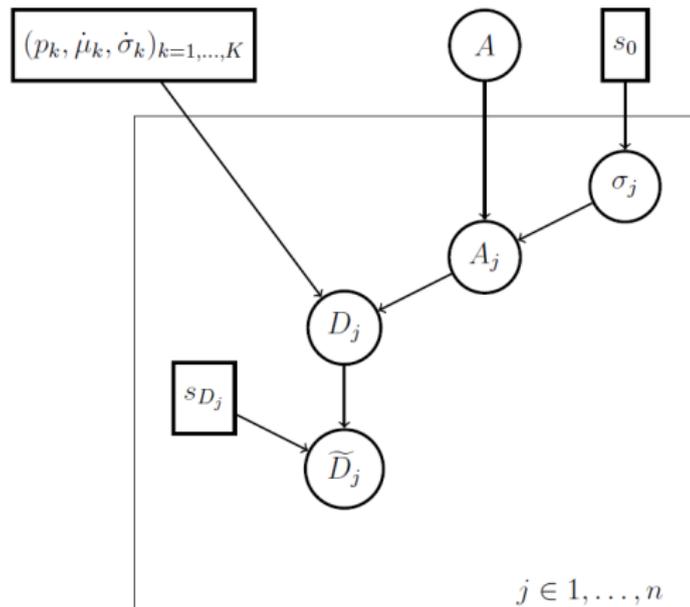
$$D_j \sim \sum_{k=1}^K \dot{p}_k \mathcal{N}(A_j \dot{\mu}_k, A_j^2 \dot{\sigma}_k^2)$$

$$A_j \sim \mathcal{N}(A, \sigma_j^2)$$

$$\sigma_j^2 \sim \mathcal{S}(s_0^2)$$

$$A \sim \text{Uniform}[\underline{A}, \bar{A}]$$

Mercier N., Galharret J.-M., Tribolo C., Kreutzer S., and Philippe A. (2022) Luminescence age calculation through Bayesian convolution of equivalent dose and dose-rate distributions : the De\_Dr model *Geochronology*, 4, 297–310, 202



## Autre modélisation : âge OSL

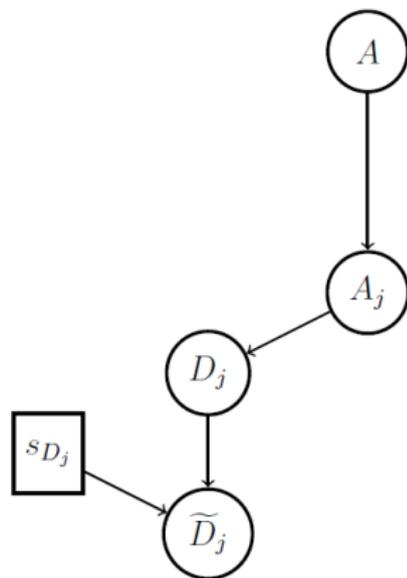
### Estimation de $s_0^2$

$$\tilde{D}_j \sim \mathcal{N}(D_j, s_{D_j}^2),$$

$$D_j \sim \sum_{k=1}^K \dot{p}_k \mathcal{N}(A_j \dot{\mu}_k, A_j^2 \dot{\sigma}_k^2),$$

$$A_j \sim \text{Uniform}[\underline{A}, \bar{A}]$$

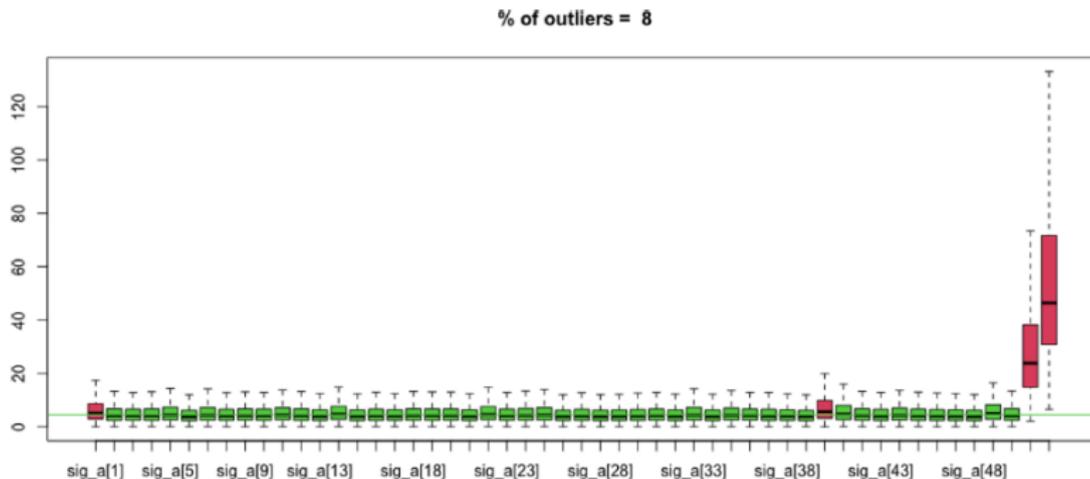
$s_0^2$  est estimée par la moyenne harmonique de  $\text{Var}(A_j \mid \tilde{D}_j)$



# Autre modélisation : âge OSL

## Estimation sur données réelles

On considère un échantillon de  $n = 53$  observations.



**Figure** – Boxplot des lois a posteriori des  $\sigma_j$  ordonnées selon les médianes a posteriori des âges.

## Validation finale du modèle

On revient à la relation fondamentale  $D \stackrel{\mathcal{L}}{=} A \dot{d}$  et on va comparer

- La fonction empirique des  $(D_j)_{j \in J}$  (notée  $F_D$ )
- La fonction de répartition de  $A \dot{d}$  (notée  $F_{A \dot{d}}$ )

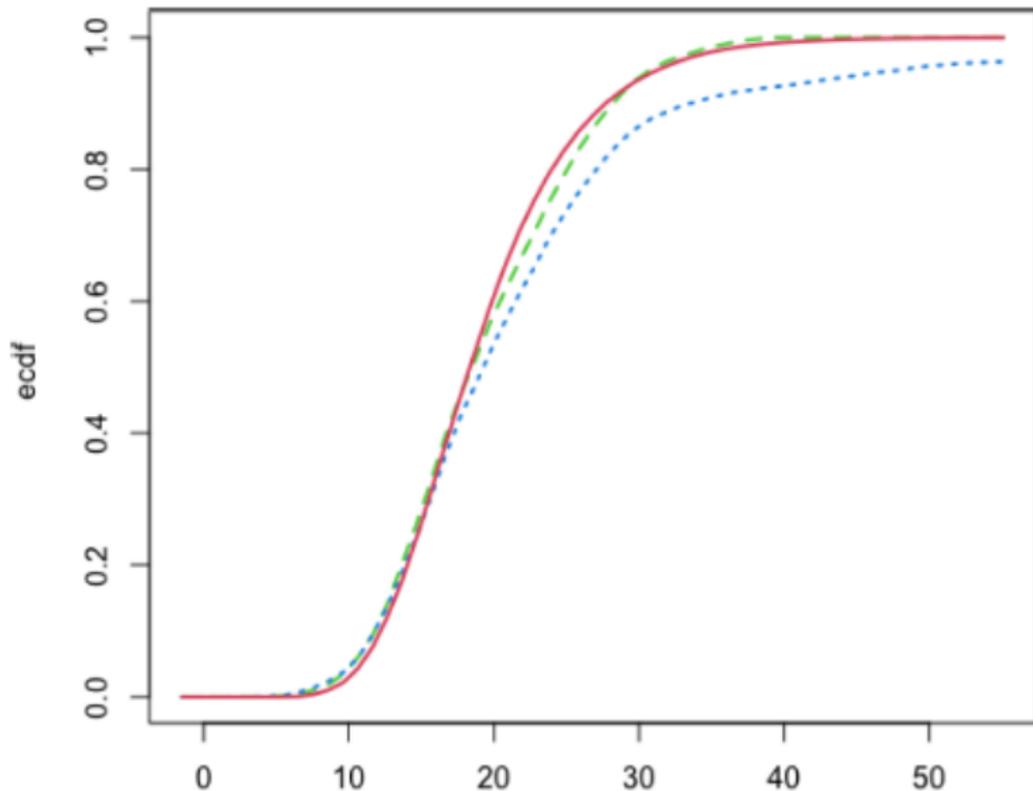
Ces deux fonctions de répartition dépendent respectivement des paramètres inconnus  $(D_j)_{j \in J}$  et  $A, \dot{\varepsilon}$ . On calcule leurs estimateurs de Bayes

$$\mathbb{E} \left( F_D | (\tilde{D}_j)_{j \in J} \right) = \frac{1}{J} \sum_{j \in J} F_{D_j | \tilde{D}_j}(t)$$

$$\mathbb{E} \left( F_{A \dot{d}} | (\tilde{D}_j)_{j \in J} \right) = \mathbb{E} \left( \sum_{k=1}^K \dot{p}_k \dot{F}_k \left( \frac{t}{A} - \varepsilon \right) | (\tilde{D}_j)_{j \in J} \right)$$

où  $\dot{F}_k$  est la fonction de répartition de  $\mathcal{N}(\dot{\mu}_k, \dot{\sigma}_k)$

**Figure** – Comparaison des fonctions de répartition empiriques de  $Ad$  (rouge),  $D$  (bleue) et  $D$  après avoir retiré les outliers détectés (vert)



- La modélisation proposée dans le cadre de la datation OSL est implémentée dans le package *Luminescence*.
- La procédure pour le contexte normal-normal a été ajouté dans le logiciel Chronomodel.
- On a aussi proposé une détection des outliers dans les modèles de comptage.
- On s'intéresse actuellement à l'adaption du modèle lorsque l'on a des contrainte stratigraphiques (contraintes d'ordre).

Merci de votre attention