

Détection bayésienne d'outliers et ses applications en archéologie

JM GALHARRET, A PHILIPPE, N MERCIER

Laboratoire Jean Leray (Univ. Nantes), CRP2A (Univ. Bordeaux)

June 4, 2019

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 mélange Bronk Ramsey (2009),
- Chronomodel : Méthode d'estimation robuste Lanos and Philippe (2017, 2018).

Notre approche :

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

Modèle hiérarchique pour estimer l'âge A d'un événement 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 A_1, \dots, A_n sont contemporains de A .

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

- Problème lié aux outliers

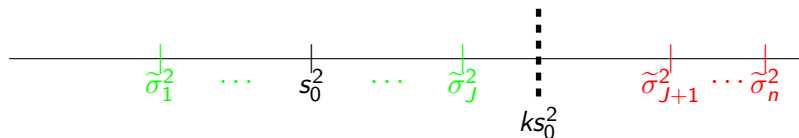


Event model :

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

où les σ_i^2 sont i.i.d de loi de shrinkage π_s telle que $\text{med}(\sigma_i^2) = s_0^2$ où s_0^2 est la moyenne harmonique des s_i^2 .

Soit $\tilde{\sigma}_i^2$ la médiane de la loi a posteriori de σ_i^2 .



Règle de décision

X_i est un outlier si : $\tilde{\sigma}_i^2 > ks_0^2$, où $k \geq 1$ est calibré pour que si l'échantillon ne contient pas d'outlier, on ne détecte qu'au plus 5% de faux positifs.

k est calibré numériquement par la méthode de Monte-Carlo :

- On simule B échantillons indépendants de taille n ne contenant aucun outlier.
- Pour tout $k \geq 1$, on calcule le pourcentage moyen $\bar{p}_{n,k}$ d'outliers détectés à tort.

$$\operatorname{argmin}_{k \geq 1} (\bar{p}_{n,k} < 0.05).$$

- On ne conserve que le sous-échantillon $(X_j)_{j \in J}$ correspondant aux mesures qui n'ont pas été détectées comme outliers.
- Sur ce sous-échantillon on estime A avec (selon l'avis de l'expert):
 - ① M2 : objets contemporains de l'événement

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

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

$$\sigma \sim \pi_\sigma$$

$$A \sim \mathcal{U}[\underline{A}, \bar{A}]$$

- ② M3 : âge identique.

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

$$A \sim \mathcal{U}[\underline{A}, \bar{A}]$$

On simule sous le modèle M_3 des échantillons de taille n et d'âge commun $A = 0$ (sans perte de généralité) contaminés par τ % d'outliers selon le modèle :

$$X_i \sim (1 - \tau)\mathcal{N}(0, s_i^2) + \tau\mathcal{N}(15, s_i^2) \quad (1)$$

On va comparer les résultats numériques des modèles :

- Event model (M_1): $\forall i = 1, \dots, n : \sigma_i^2$ individuels.
- M_2 : $\forall i \in J : \sigma_i^2 = \sigma^2$
- M_3 : $\forall i \in J : \sigma_i^2 = 0$

Calibration du seuil k

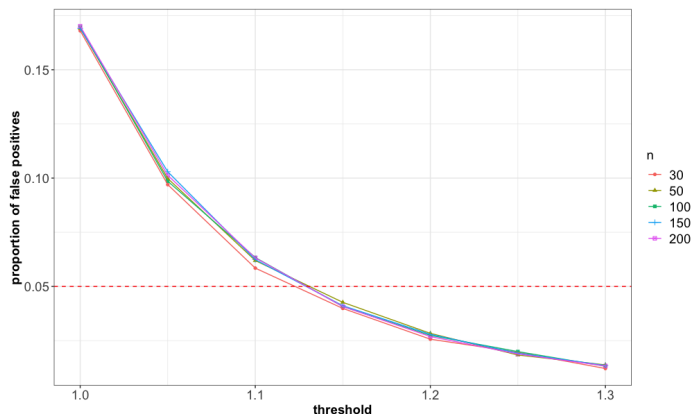


Figure: % moyen de faux positifs $\bar{p}_{n,k} = \sum_{b=1}^B p_{n,k}^{(b)} / B$ détectés pour $B = 1000$ répliques indépendantes sous \mathcal{H}_0 en fonction de la taille n de l'échantillon et de $k \geq 1$.

Comparaison des approches pour l'estimation de l'âge

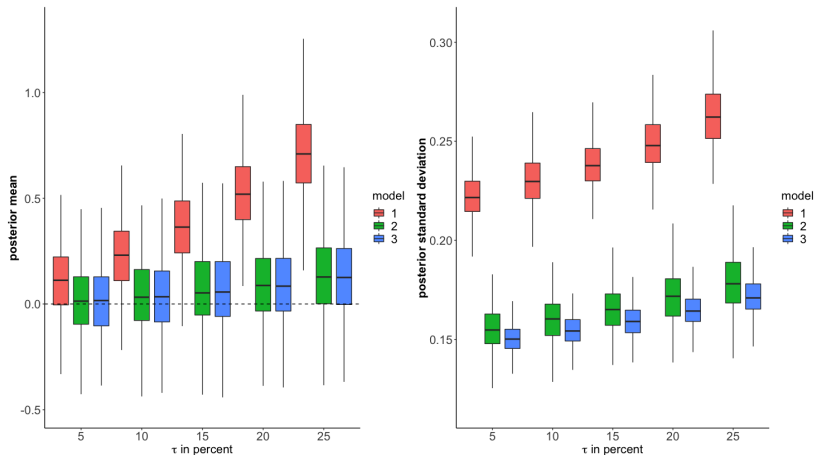


Figure: Boxplot des moyennes (gauche) et des écarts types (droite) a posteriori en fonction du taux de contamination τ .

Application à la datation par luminescence

Relation fondamentale

$$D \stackrel{\mathcal{L}}{=} Ad \quad (2)$$

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

- Le débit de dose \dot{d} n'est pas observable,
- On sait simuler des échantillons \dot{d} avec une erreur systématique $\dot{\varepsilon}$.
- Modélisation de \dot{d} :

$$\tilde{\dot{d}} = \dot{d} + \dot{\varepsilon} \sim \mathcal{C}(\dot{\mu} + \dot{\varepsilon}, \dot{\sigma}^2), \quad \dot{\varepsilon} \sim \mathcal{N}(0.10\dot{\mu}, 0.01\dot{\mu}^2)$$

- La relation (2) devient

$$D \sim \mathcal{C}(A(\dot{\mu} - \dot{\varepsilon}), A^2\dot{\sigma}^2)$$

Modèle d'âge

Les doses absorbées $(D_j)_{j \in \{1, \dots, n\}}$ sont mesurées avec une erreur gaussienne de variance connue :

$$\tilde{D}_j \sim \mathcal{N}(D_j, s_j^2), j \in \{1, \dots, n\}$$

où \tilde{D}_j désigne la dose absorbée mesurée.

$$D_j \sim \mathcal{C}(A_j(\mu - \varepsilon), A_j^2 \sigma^2)$$

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

$$\sigma_j^2 \sim \mathcal{S}(3, S_0^2)$$

$$A \sim \mathcal{U}[A, \bar{A}]$$

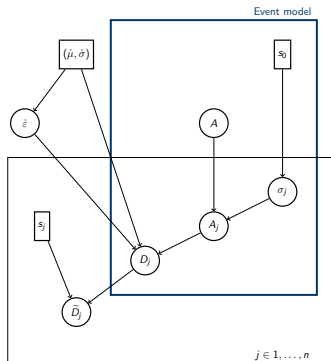


Figure: DAG of the model

Calibration du seuil

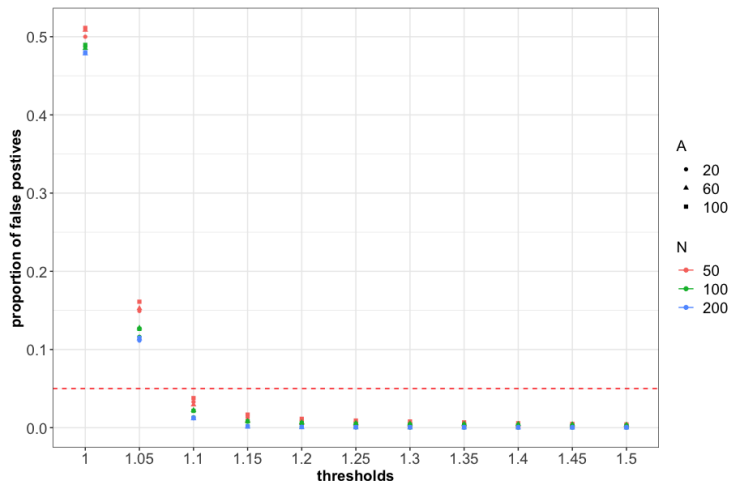
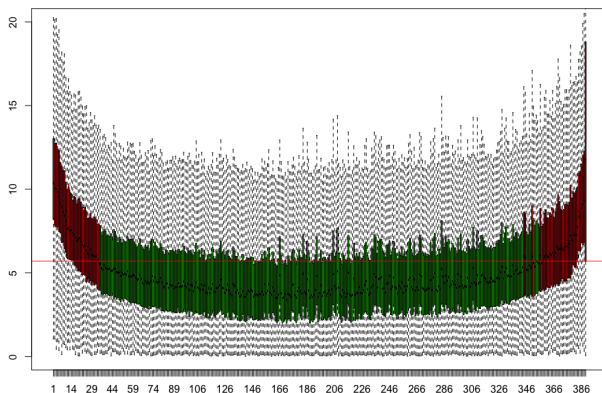


Figure: Calibration du seuil k en fonction de A , n .

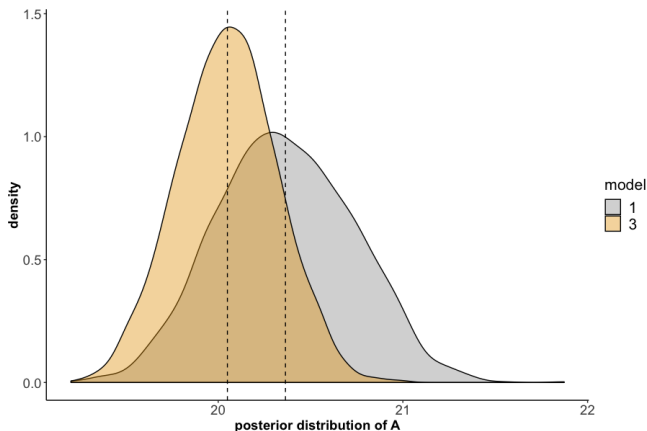
Analyse de données

On considère un échantillon de $n = 389$ doses absorbées. Boxplot des distributions a posteriori de σ_j ordonnées selon la médiane a posteriori des âges A_j . La couleur rouge indique les outliers détectés (18 %). La ligne rouge correspond à $1.1s_0$.



Comparaison des estimations

Comparaison des lois a posteriori de A : en or l'estimation sans les outliers ($\hat{A} = 20.08$, 95%-CI = [19.58, 20.64]) et en gris celle avec le modèle robuste 20.54 and $\hat{A} = 20.54$, 95%-CI = [19.77, 21.34])



Pour valider le choix du modèle on va comparer :

- La fonction de répartition empirique de $(D_j)_{j \in J}$ (notée F_D)
- La fonction de répartition de Ad (notée F_{Ad}).

F_D et F_{Ad} dépendent des paramètres inconnus $(D_j)_{j \in J}$ et $(A, \dot{\epsilon})$.

On calcule leur estimateur de Bayes:

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

$$\mathbb{E}(F_{Ad}(t) | (\tilde{D}_j)_{j \in J}) = \mathbb{E} \left(\dot{G} \left(\frac{t}{A} - \dot{\epsilon} \right) | (\tilde{D}_j)_{j \in J} \right)$$

où \dot{G} est la fonction de répartition de la loi de Cauchy $(\dot{\mu}, \dot{\sigma})$.

Validation du modèle

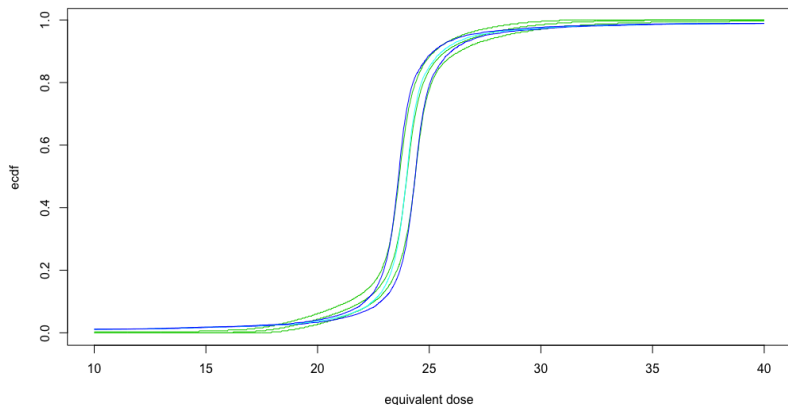


Figure: Représentation des intervalles de crédibilité à 95% de F_D , F_{Ad} , respectivement en bleu et en vert.

- Bronk Ramsey, C. (2009). Dealing with outliers and offsets in radiocarbon dating. *Radiocarbon*, 51(3):1023–1045.
- Lanos, P. and Philippe, A. (2017). Hierarchical Bayesian modeling for combining dates in archaeological context. *Journal de la Societe Française de Statistique*, 158(2):72–88.
- Lanos, P. and Philippe, A. (2018). Event date model: a robust bayesian tool for chronology building. *Communications for Statistical Applications and Methods*, 25(2):131–157.