

Données manquantes et régression multinomiale

Analyse multidimensionnelle appliquée

Léo Belzile

HEC Montréal

- Les données manquantes réduisent la quantité d'information disponible et augmentent l'incertitude.
- On ne peut pas les ignorer (étude des cas complets) sans biaiser les interprétations et réduire la quantité d'information disponible.
- On considère des méthodes d'imputation qui remplacent les valeurs manquantes.

Considérons le cas d'une régression logistique pour une variable explicative binaire.

Plutôt que d'assigner à la classe la plus probable, une prédiction aléatoire simule une variable 0/1 avec probabilité $(1 - \hat{p}_i, \hat{p}_i)$.

```
1 pred <- 0.3 #probabilité de succès  
2 rbinom(n = 15, size = 1, prob = pred)
```

```
[1] 0 1 1 0 0 0 0 1 1 1 1 0 1 0 1
```

On ne tient pas compte du fait que des valeurs ont été remplacées (on fait comme si c'était de vraies observations).

On sous-évalue encore une fois la variabilité des données

- les écarts-type des estimations sont trop petits.

Inspection des valeurs manquantes

Il est donc nécessaire d'examiner la configuration des valeurs manquantes avant de faire quoi que ce soit.

```
1 data(manquantes, package = 'hecmulti')
2 summary(manquantes)
3 # Pourcentage de valeurs manquantes
4 apply(manquantes, 2, function(x){mean(is.na(x))})
5 # Voir les configurations de valeurs manquantes
6 md.pattern(manquantes) # graphique diapo suivante
```

Table 1: Nombre et pourcentage de valeurs manquantes par variable.

| | x1 | x2 | x3 | x4 | x5 | x6 | y |
|-------------|------|-----|----|------|----|----|---|
| nombre | 192 | 49 | 0 | 184 | 0 | 0 | 0 |
| pourcentage | 38.4 | 9.8 | 0 | 36.8 | 0 | 0 | 0 |

Configuration des valeurs manquantes

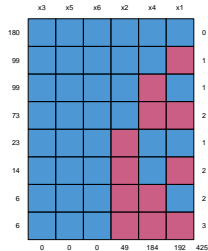
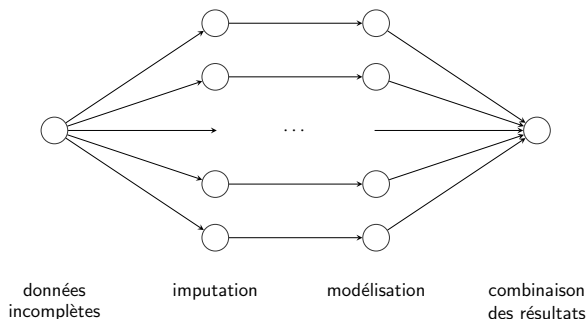


Figure 1: Configurations des valeurs manquantes pour manquantes.

Imputation multiple

Valides pour les données manquantes de manière aléatoire et complètement aléatoires (MAR et MCAR).

1. Procéder à plusieurs imputations aléatoires pour obtenir un échantillon complet (*mice*)
2. Ajuster le modèle d'intérêt avec chaque échantillon (*with*). 3. Combiner les résultats obtenus (*pool* et *summary*)



Considérons un seul paramètre θ (ex: coefficient d'une régression) et supposons qu'on procède à K imputations.

On estime les paramètres du modèle séparément pour chacun des K ensembles de données imputés, disons

- $\hat{\theta}_k$ pour l'estimation du paramètre θ dans l'échantillon k et
- $\hat{\sigma}_k^2 = \text{Va}(\hat{\theta}_k)$ pour l'estimation de la variance de $\hat{\theta}_k$.

L'estimation finale de θ , dénotée $\hat{\theta}$, est obtenue tout simplement en faisant la moyenne des estimations de tous les modèles, c'est-à-dire,

$$\hat{\theta} = \frac{\hat{\theta}_1 + \dots + \hat{\theta}_K}{K}.$$

Une estimation ajustée de la variance de $\hat{\theta}$ est

$$\text{Va}(\hat{\theta}) = W + \frac{K+1}{K}B,$$

où

$$W = \frac{1}{K} \sum_{k=1}^K \hat{\sigma}_k^2 = \frac{\hat{\sigma}_1^2 + \dots + \hat{\sigma}_K^2}{K},$$

$$B = \frac{1}{K-1} \sum_{k=1}^K (\hat{\theta}_k - \hat{\theta})^2.$$

- W est la moyenne des variances (variance intra-groupe) et
- B la variance des moyennes (variance inter-groupe).

Des formules analogues existent pour les degrés de liberté, les valeurs- p , etc. ainsi que pour la cas multidimensionnel (plusieurs paramètres).

Si on procédait à une seule imputation (même en ajoutant une part d'aléatoire pour essayer de reproduire la variabilité des données), on ne serait pas en mesure d'estimer la variance inter-groupe de l'estimateur.

On peut estimer la fraction de l'information manquante sur θ avec $(1 + 1/K)B/Va(\hat{\theta})$.

Avec p variables X_1, \dots, X_p , spécifier un ensemble de modèles conditionnels pour chaque variable X_j en fonction de

- toutes les autres variables, X_{-j}
 - les valeurs observées pour cette variable, $X_{j,\text{obs}}$
1. Initialisation: remplir les trous avec des données au hasard parmi $X_{j,\text{obs}}$ pour $X_{j,\text{man}}$
 2. À l'itération t , pour chaque variable $j = 1, \dots, p$, à tour de rôle:
 - a) tirage aléatoire des paramètres $\phi_j^{(t)}$ du modèle pour $X_{j,\text{man}}$ conditionnel à $X_{-j}^{(t-1)}$ et $X_{j,\text{obs}}$
 - b) échantillonnage de nouvelles observations $X_{j,\text{man}}^{(t)}$ du modèle avec paramètres $\phi_j^{(t)}$ conditionnel à $X_{-j}^{(t-1)}$ et $X_{j,\text{obs}}$
 3. Répéter le cycle

Imputation multiple avec mice

```
1 library(mice)
2 # Intensif en calcul, réduire "m" si nécessaire
3 impdata <- mice(
4   data = manquantes,
5   # argument "method" pour le modèle
6   # dépend du type des variables, par ex.
7   # régression logistique pour données binaires
8   m = 50, # nombre d'imputations
9   seed = 60602, # germe aléatoire
10  printFlag = FALSE)
11 # Extraite une copie (m=1,..., 50) imputée
12 complete(data = impdata,
13           action = 1) #no de la copie
```

```
1 # ajuste le modèle avec les données imputées
2 adj_im <- with(
3   data = impdata,
4   expr = glm(y ~ x1 + x2 + x3 + x4 + x5 + x6,
5             family = binomial))
6 # combinaison des résultats
7 fit <- pool(adj_im)
8 summary(fit)
```

| terme | estimation | erreur-type | stat | ddl | valeur-p |
|-------------|------------|-------------|-------|--------|----------|
| (Intercept) | -2.86 | 1.08 | -2.65 | 327.28 | 0.009 |
| x12 | -0.90 | 0.58 | -1.55 | 172.63 | 0.124 |
| x13 | -0.53 | 0.59 | -0.90 | 164.98 | 0.372 |
| x14 | -0.76 | 0.55 | -1.39 | 241.12 | 0.166 |
| x15 | -0.58 | 0.64 | -0.91 | 176.76 | 0.365 |
| x22 | -0.03 | 0.44 | -0.07 | 403.25 | 0.944 |
| x23 | -0.65 | 0.49 | -1.34 | 353.29 | 0.181 |
| x24 | -2.60 | 0.65 | -4.00 | 271.66 | <1e-04 |
| x25 | -1.32 | 0.72 | -1.84 | 283.49 | 0.066 |
| x3 | 1.61 | 0.31 | 5.19 | 300.31 | <1e-04 |
| x4 | 2.56 | 0.42 | 6.14 | 185.05 | <1e-04 |
| x5 | 0.11 | 0.02 | 5.27 | 436.59 | <1e-04 |
| x62 | -1.41 | 0.37 | -3.83 | 408.17 | 1e-04 |
| x63 | -2.58 | 0.41 | -6.24 | 399.50 | <1e-04 |

- Les données manquantes réduisent la quantité d'information disponible et augmentent l'incertitude.
- On ne peut pas les ignorer (étude des cas complets) sans biaiser les interprétations et réduire la quantité d'information disponible.
- Pour bien capturer l'incertitude et ne pas modifier les relations entre variables, il faut utiliser une méthode aléatoire.
- Avec l'algorithme MICE, on utilise un modèle conditionnel pour chaque variable à tour de rôle

L'imputation multiple est préférée à l'imputation simple car elle permet d'estimer l'incertitude sous-jacente en raison des données manquantes.

- On procède à l'imputation plusieurs fois (avec un modèle conditionnel, prédictions différentes chaque fois)
- on crée plusieurs copies
- ajuste le modèle sur chacune et
- combine les résultats

Traitement spécial pour erreurs-type, degrés de liberté, valeurs- p et intervalles de confiance.

On considère une variable réponse catégorielle avec $K \geq 2$ modalités.

Objectif: modéliser la probabilité de chaque catégorie de la variable réponse.

Soit la probabilité d'appartenir à la modalité k ,

$$p_{ik} = \Pr(Y_i = k \mid X_i), \quad (k = 1, \dots, K).$$

La somme des probabilités, $p_{i0} + \dots + p_{iK}$, vaut 1.

Comme avec la régression logistique, on fixe une catégorie de référence (disons 1) et on modélise le log de la cote de chacune des autres catégories par rapport à cette référence,

$$\ln \left(\frac{p_{ij}}{p_{i1}} \right) = \eta_{ij} = \beta_{0j} + \beta_{1j}x_{i1} + \cdots + \beta_{pj}x_{ip}, \quad (j = 2, \dots, K).$$

- Avec K modalités et p variables explicatives, on obtiendra $(K - 1) \times (p + 1)$ paramètres à estimer, en incluant l'ordonnée à l'origine.

L'interprétation des paramètres se fait comme en régression logistique sauf qu'il faut y aller équation par équation.

On peut aussi exprimer le modèle en termes des probabilités,

$$p_{i1} = \Pr(Y_i = 1 \mid X_i) = \frac{1}{1 + \exp(\eta_{i2}) + \cdots + \exp(\eta_{iK})}$$

$$p_{ik} = \Pr(Y_i = k \mid X_i) = \frac{\exp(\eta_{ik})}{1 + \exp(\eta_{i2}) + \cdots + \exp(\eta_{iK})}, \quad k = 2, \dots, K.$$

où η_{ij} est le prédicteur linéaire de l'individu i pour le log de la cote de $Y_i = j$ versus la référence $Y_i = 1$.

Les données de cet exemple sont tirées d'un sondage Ipsos réalisé pour le site de nouvelles FiveThirtyEight.

La base de données `vote` contient 5837 observations avec les pondérations associées.

Nous allons modéliser l'intention de vote, `catvote` à l'aide d'une régression logistique multinomiale.

Analyse exploratoire

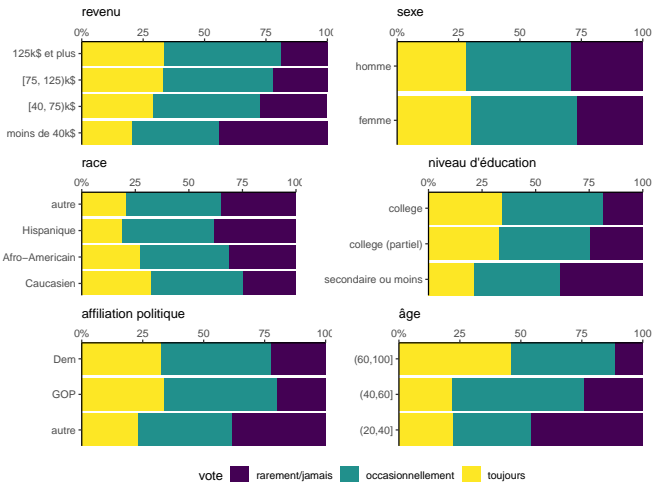


Figure 2: Proportion des modalités des variables sociodémographiques des données de participation électorale.

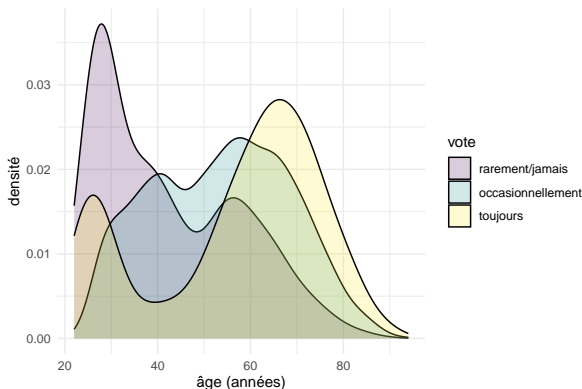


Figure 3: Fréquence de vote selon l'âge.

Notez le comportement des jeunes voteurs (bimodal). Ces personnes n'ont souvent eu qu'une seule occasion de voter...

Ajustement du modèle

La fonction `multinom` du paquet `nnet` ajuste le modèle multinomial logistique.

```
1 data(vote, package = "hecmulti")
2 levels(vote$catvote)
```

```
[1] "rarement/jamais" "occasionnellement" "toujours"
```

```
1 # Modèle multinomial
2 multi1 <- nnet::multinom(
3   catvote ~ age + sexe + race + revenu +
4   educ + affiliation,
5   data = vote,           # base de données
6   subset = age > 30,     # sous-ensemble des données
7   weights = poids,       # poids de sondage
8   trace = FALSE)         # infos sur convergence
```



```
1 # Tableau résumé de l'ajustement
2 summary(multi1)
3 # Estimations des coefficients
4 coef(multi1)
5 # Intervalles de confiance (Wald)
6 confint(multi1)
7 # Critères d'information
8 AIC(multi1)
9 BIC(multi1)
10 # Prédiction: probabilité de chaque modalité
11 predict(multi1, type = "probs")
12 # Prédiction: classe la plus susceptible
13 predict(multi1, type = "class")
```

Comparaison de modèles emboîtés

Le modèle avec uniquement l'ordonnée à l'origine possède $K - 1$ paramètres. Il retourne comme probabilité prédite la proportion empirique de chaque catégorie.

```
1 multi0 <- nnet::multinom(catvote ~ 1,  
2                               weights = poids,  
3                               subset = age > 30,  
4                               data = vote,  
5                               trace = FALSE)  
6 # Test de rapport de vraisemblance  
7 anova(multi0, multi1)
```

| Resid. df | Resid. Dev | Test | Df | LR stat. | Pr(Chi) |
|-----------|------------|--------|----|----------|---------|
| 9692 | 9781.07 | | | | |
| 9668 | 8504.74 | 1 vs 2 | 24 | 1276.33 | 0 |

Pour un profil X_i donné, on peut

- calculer chacun des $K - 1$ prédicteurs linéaires $\hat{\eta}_{i2}, \dots, \hat{\eta}_{iK}$.
- écrire $p_{ik} = p_{i1} \exp(\hat{\eta}_{ik})$ (formule de la cote)
- substituer cette mesure dans l'équation $p_{i1} + \dots + p_{iK} = 1$
- isoler la prédiction numérique pour p_{i1} .
- en déduire les probabilités de succès de chaque modalité de Y .

Exemple au tableau

La prédiction du modèle est une probabilité pour chacune des K modalités.

On peut toujours classifier les événements

- avec $K - 1$ points de coupure...
- ou assigner à la modalité la plus probable

Avec les prédictions, on peut comparer les observations et les prédictions à l'aide d'une matrice de confusion $K \times K$.

- Le taux de bonne classification est toujours valide
- Il existe des extensions multidimensionnelles de l'aire sous la courbe

- Contrairement à la régression logistique, le nombre de paramètres augmente rapidement avec le nombre de variables explicatives, p .
- Il y a moins d'information pour estimer les paramètres qu'une régression linéaire: prévoir de plus grandes tailles d'échantillon.
- Attention aux modalités à faible fréquence et à la répartition des variables explicatives au sein des différentes modalités.

Outre la régression multinomiale logistique, on peut également considérer la régression logistique cumulative à cotes proportionnelles.

- modèle plus parcimonieux que le modèle multinomial logistique,
- mais au prix de postulats supplémentaires...

En R, la variable réponse doit être de classe `ordered`, un facteur dont les niveaux sont ordonnés en ordre croissant.

```
1 class(hecmulti::vote$catvote)
```

```
[1] "ordered" "factor"
```

- Soit p_1, \dots, p_K les probabilités associées aux événements $Y = 1, \dots$
- On définit les points de coupure pour les K classes,

$$-\infty = \zeta_0 < \zeta_1 < \dots < \zeta_K = \infty.$$

- Il y a $K - 1$ paramètres ζ à déduire pour identifier les probabilités puisque $p_1 + \dots + p_K = 1$.

Le modèle logistique à cote proportionnelle spécifie $K - 1$ équations logistiques; pour $k = 1, \dots, K - 1$,

$$\ln \left(\frac{\Pr(Y_i > k \mid X_i)}{\Pr(Y_i \leq k \mid X_i)} \right) = -\zeta_k + \beta_1 X_{i1} + \dots + \beta_p X_{ip}.$$

- Les paramètres associés aux variables explicatives, β_1, \dots, β_p sont les mêmes pour chacune des log-cotes
- mais il y a une ordonnée à l'origine différente par rapport de cote, $-\zeta_k$.

On considère la cote de $\Pr(Y_i > k \mid X_i)$ versus $\Pr(Y_i \leq k \mid X_i)$, qui mesure à quel point il est plus probable que Y_i prenne une valeur supérieure à k par rapport à une valeur inférieure ou égale à k , avec

Pour chaque augmentation d'une unité de X_j , cette cote est multipliée par $\exp(\beta_j)$, peu importe la valeur de Y (cote proportionnelle).

Table 2: Tableau des estimations des coefficients du modèle pour réponses ordinales pour la régression logistique à cotes proportionnelles avec sexe.

| effet | coefficient | erreur-type |
|---|-------------|-------------|
| sexe [homme] | -0.166 | 0.055 |
| cst [rarement jamais occasionnellement] | -1.297 | 0.044 |
| cst [occasionnellement toujours] | 0.865 | 0.044 |

Les hommes de plus de 30 ans sont moins susceptibles de voter fréquemment que les femmes.

La cote catégorie plus fréquente de vote (vs moins fréquente) pour les hommes est $\exp(-0.166) = 0.847$ fois celle des femmes, soit une diminution de 15.3% de la cote.

Pour simplifier, on utilise uniquement `sexe` comme variable explicative.

```
1 # with(vote, is.ordered(catvote))
2 multi2a <- MASS::polr(
3   catvote ~ sexe,
4   data = vote,
5   subset = age > 30,
6   weights = poids,
7   method = "logistic",
8   Hess = TRUE)
9 summary(multi2a)
```

```
1 # IC pour beta_x (vraisemblance profilée)
2 confint(multi2a)
3 # On peut obtenir les intervalles de Wald
4 # avec confint.default (PAS RECOMMANDÉ)
5
6 # Critères d'information
7 AIC(multi2a); BIC(multi2a)
8 # Tableau des coefficients
9 # Coefficients (variables explicatives)
10 coef(multi2a)
11 # Négatif de l'ordonnée à l'origine:
12 multi2a$zeta
```

Si on écrit les équations pour la cote, on obtient

$$\frac{\Pr(Y = \text{rarement} \mid \text{sexe})}{\Pr(Y \geq \text{occasionnellement} \mid \text{sexe})} = \exp(-0.166\text{sexe} + 1.297)$$
$$\frac{\Pr(Y \leq \text{occasionnellement} \mid \text{sexe})}{\Pr(Y = \text{toujours} \mid \text{sexe})} = \exp(-0.166\text{sexe} - 0.865).$$

En terme de probabilité cumulée d'excéder k ,

$$\Pr(Y_i > k \mid \mathbf{X}_i) = \text{expit}(-\eta_k + \beta_1 X_{i1} + \dots + \beta_p X_{ip}), \quad k = 1, \dots, K - 1.$$

En utilisant ces expressions, on peut obtenir la probabilité de chaque catégorie,

$$\Pr(Y_i = k \mid \mathbf{X}_i) = \Pr(Y_i > k \mid \mathbf{X}_i) - \Pr(Y_i > k - 1 \mid \mathbf{X}_i).$$

- Soit $p_1 = \Pr(Y = \text{rarement/jamais} \mid \text{femme})$, etc.
- On a $\text{expit}(\zeta_k)$ ($k = 0, \dots, K$) qui donne 0, 0.215, 0.704 et 1.
- Les différences donnent $\hat{p}_1 = 0.215$, $\hat{p}_2 = 0.489$ et $\hat{p}_3 = 0.296$.
- Un rapide calcul numérique montre que c'est bien ce que retourne les prédictions.

```
1 predict(multi2a,  
2         newdata = data.frame(sexe = factor("femme")),  
3         type = "probs")
```

Une des hypothèses de ce modèle est que les effets des variables explicatives sont les mêmes pour chaque équation.

- \mathcal{H}_0 : l'effet de chaque variable est le même pour les K logit du modèle .

Une très petite valeur- p (rejet de \mathcal{H}_0) pour ce test serait une indication que le modèle multinomial logistique serait préférable.

Test de rapport de vraisemblance

Ce test compare les deux modèles emboîtés, avec

- hypothèse nulle \mathcal{H}_0 : modèle cumulatif à cotes proportionnelles, avec $p + K - 1$ paramètres
- hypothèse alternative \mathcal{H}_a : modèle multinomial, avec $(K - 1) \times (p + 1)$ paramètres

```
1 multi2b <- nnet::multinom(catvote ~ sexe,  
2   data = vote, subset = age > 30,  
3   weights = poids, trace = FALSE)  
4 # Valeur-p du test de rapport de vraisemblance  
5 pchisq(q = deviance(multi2a) - deviance(multi2b),  
6   df = length(coef(multi2a)),  
7   lower.tail = FALSE)
```

```
[1] 0.7886889
```

Le modèle sous \mathcal{H}_0 semble être une simplification adéquate.

Si on ajuste plutôt le modèle avec uniquement *age*, la valeur-*p* est inférieure à 10^{-5} : le modèle cumulatif à cote proportionnelles ne serait pas une simplification adéquate.

On peut également effectuer des tests pour déterminer la significativité

- la significativité globale (ordonnée à l'origine vs modèle complet)
- l'effet d'une variable explicative (modèle complet, moins une variable)

Comparaison des prédictions

Prédictions pour le modèle avec uniquement age comme variable explicative.

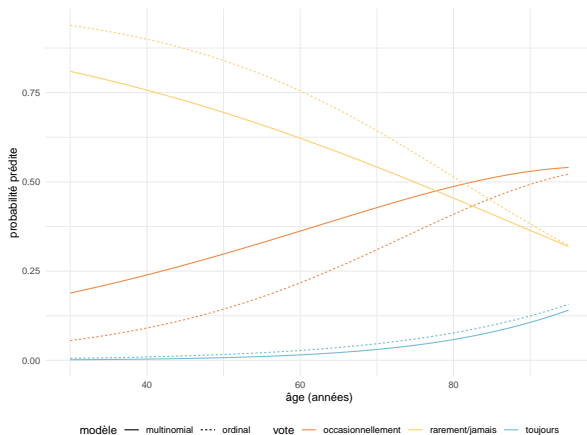


Figure 4: Probabilités prédites pour chaque modalité selon l'âge.

- La régression multinomiale logistique pour une variable catégorielle à K niveaux est une extension directe de la régression logistique pour données binaires
 - la somme des probabilités vaut 1.
 - il y a $K - 1$ équations de cote en termes des variables explicatives,
 - donc le nombre de paramètres croît rapidement.

On met beaucoup l'accent sur l'interprétation des coefficients à l'échelle de la cote.

- rapports de cote = modèle multiplicatif: la cote de catégorie k vs référence est multipliée par $\exp(\beta_{jk})$ pour chaque augmentation de X_j d'une unité.
- les coefficients manquants (cote de $Y = k$ vs $Y = l$) peut être déduits par des manipulations algébriques.

Les outils usuels d'inférence pour les modèles estimés par maximum de vraisemblance sont applicables.

- intervalles de confiance (Wald ou vraisemblance profilée)
- tests de rapport de vraisemblance
- critères d'information

Côté classification, on va règle générale assigner à la classe la plus probable.

- il existe des équivalents multidimensionnels directs à ce qu'on a couvert (matrice de confusion, taux de bonne classification, gain, etc.)
- certains concepts (sensibilité, spécificité, fonction d'efficacité du récepteur) ne sont en revanche pas applicables ou n'ont pas d'équivalent.

Le modèle cumulatif à cote proportionnelle est une simplification du modèle multinomial pour des données ordinales.

- On suppose que l'effet des variables est le même pour la cote de la survie de chaque modalité.
- Moins de paramètres, mais postulat à vérifier (via test de rapport de vraisemblance).