

# Loi des EMV

Marie-Anne Poursat

M1 Bioinformatique [Université Paris-Saclay](#)

30 janvier 2024

# Estimateur du maximum de vraisemblance

Modèle *non-linéaire* en les paramètres : *pas de calculs exacts*.

↪ l'étude *asymptotique* ( $n$  "grand") est la référence

↪ les calculs sont des *approximations asymptotiques* ( $n \rightarrow \infty$ )

Modèle :  $(Y_1, \dots, Y_n)$  i.i.d. de densité  $f_\theta$ ,  $\theta = (\theta_1, \dots, \theta_p)^T$

L'EMV  $\hat{\theta}$  maximise  $\log L(\theta) = \sum_{i=1}^n \log f_\theta(Y_i)$ , solution des équations de vraisemblance :

$$G(\hat{\theta}) = \frac{\partial \log L}{\partial \theta}(\hat{\theta}) = 0$$

- Consistance (garantie d'un bon estimateur).
  - Approximation de la loi de  $\hat{\theta}$ ?
    - ↪ indispensable pour calculer des intervalles de confiance ou des tests.
- Comment calculer la variance de  $\hat{\theta}$ ?

# Information de Fisher

On appelle **information de Fisher** de l'échantillon  $(Y_1, \dots, Y_n)$  la quantité

$$I_n(\theta) = \text{Var} \left[ \frac{\partial \log L}{\partial \theta}(\theta) \right] = \text{E} \left[ \left( \frac{\partial \log L}{\partial \theta}(\theta) \right)^2 \right]$$

- si  $\theta$  est un réel,  $I(\theta)$  est un nombre ;
- si  $\theta$  est un vecteur,  $I(\theta)$  est une matrice  $p \times p$ , définie positive (en  $\hat{\theta}$ ).

On peut montrer :

$$I(\theta) = -\text{E} [H(\theta)] = -\text{E} \left[ \frac{\partial^2}{\partial \theta_j^2} \log L(\theta) \right]_j$$

*Interprétation géométrique* : plus  $I(\theta)$  est grande, meilleure est la localisation du maximum de la log-vraisemblance.

# Calcul de $I$

Comment calculer  $I$ ? En général, on ne sait pas calculer l'espérance (dépend de  $\theta$  inconnu, calcul analytique impossible), on l'estime par l'information de Fisher observée

$$\hat{I} = -H(\hat{\theta})$$

calculée au dernier pas de l'algorithme d'optimisation.

Exemple : échantillon de loi de Bernoulli.

## Théorème ((généralisation du TLC))

*Sous des hypothèses mathématiques de régularité du modèle (log L 2 fois dérivable, le support de la loi ne dépend pas de  $\theta$ ,  $0 < I(\theta) < \infty$ ), dans le cas où  $\theta$  est un réel,*

❶ *L'EMV  $\hat{\theta}$  est un estimateur consistant de  $\theta$ ,*

❷  $\widehat{\text{s.e.}} = \text{s.e.}(\hat{\theta}) \approx \sqrt{\frac{1}{\hat{I}}}$

❸  $\frac{\hat{\theta} - \theta}{\widehat{\text{s.e.}}} \overset{\text{Loi}}{\rightsquigarrow} \mathcal{N}(0, 1)$

*Ce qui se traduit par une approximation de la loi de  $\hat{\theta}$  :*

$$\text{Loi}(\hat{\theta}) \approx \mathcal{N}\left(\theta, \frac{1}{\hat{I}}\right)$$

Les logiciels statistiques implémentent le calcul de l'écart-type (estimé) de l'EMV

On déduit de l'approximation de la loi de l'EMV l'intervalle de confiance de niveau (approché)  $1 - \alpha$  suivant :

$$\left( \hat{\theta} - q_{1-\alpha/2} \widehat{\text{s.e.}}, \quad \hat{\theta} + q_{1-\alpha/2} \widehat{\text{s.e.}} \right)$$

avec  $q_{1-\alpha/2}$  le quantile d'ordre  $(1 - \alpha/2)$  d'une loi  $\mathcal{N}(0, 1)$ .

$$\theta = (\theta_1, \dots, \theta_p)^T$$

$$I(\theta) = \begin{pmatrix} -E(H_{11}(\theta)) & -E(H_{12}(\theta)) & \dots & -E(H_{1p}(\theta)) \\ -E(H_{12}(\theta)) & -E(H_{22}(\theta)) & \dots & -E(H_{2p}(\theta)) \\ \vdots & \vdots & & \vdots \\ -E(H_{p1}(\theta)) & -E(H_{p2}(\theta)) & \dots & -E(H_{pp}(\theta)) \end{pmatrix}$$

- $I(\theta)$  estimée par  $\hat{I}$ .
- La variance de  $\hat{\theta}$  est estimée par  $\hat{I}^{-1}$ , matrice  $p \times p$

## Loi de l'EMV

Si  $\theta \in \mathbb{R}^p$ , on peut montrer que dans les modèles réguliers,

$$\hat{\theta}^{ML} \sim \mathcal{N}(\theta, I^{-1}(\theta))$$

Cette approximation est encore valide si  $I(\theta)$  est estimée par  $\hat{I}$ .

On peut donc calculer

$$\widehat{\text{s.e.}}(\hat{\theta}_j^{ML}) \approx \sqrt{\hat{I}_{jj}}$$

et

$$IC(\theta_j) = \left[ \hat{\theta}_j^{ML} - q_{1-\alpha/2} \widehat{\text{s.e.}}(\hat{\theta}_j^{ML}); \quad \hat{\theta}_j^{ML} + q_{1-\alpha/2} \widehat{\text{s.e.}}(\hat{\theta}_j^{ML}) \right]$$

est un intervalle de confiance de niveau  $1 - \alpha$ .



# Exemple : données

Deux séquences ADN alignées de longueur  $N$

↪ *match* = positions qui présentent le même nucléotide A, C, G ou T (repérées par \*)

*		*			*			*	*	*		*	*								*		
G	G	A	G	A	C	T	G	T	A	G	A	C	A	G	C	T	A	A	T	.....	G	A	G
G	A	A	C	G	C	C	C	T	A	G	C	C	A	C	G	A	G	C	C	.....	G	G	C
1	2	3	4	5	6	7	8	9	.	.	.	.	.	.	.	.	.	.	.	.	.	.	N

$Y$  = nombre de *matches* consécutifs observés avant une position où les nucléotides diffèrent

- Observation :  $n \leq N$  réalisations  $Y_1, \dots, Y_n$  de  $Y$  le long de l'alignement
- $n$  est donc le nombre de positions de l'alignement qui ne sont pas des *match*
- Sur l'exemple :  $y_1 = 1, y_2 = 1, y_3 = 0, y_4 = 1, y_5 = 0, y_6 = 3$ , etc...

# Exemple : estimation

On suppose que  $Y_1, \dots, Y_n$  forme un échantillon indépendant et de même loi géométrique :

$$P(Y = y) = (1 - \theta)\theta^y, \quad y = 0, 1, 2, \dots$$

On a alors  $E(Y) = \frac{\theta}{1 - \theta}$ .

Paramètre :  $\theta$  est la probabilité d'observer un *match*.

- 1 Déterminer la log-vraisemblance des observations.
- 2 Montrer que l'estimateur du maximum de vraisemblance de  $\theta$  est  $\hat{\theta} = \frac{\bar{Y}}{1 + \bar{Y}}$ , où  $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$ .
- 3 Calculer l'information de Fisher des observations. En déduire une approximation de la loi de  $\hat{\theta}$ .
- 4 Donner le code R qui permet de calculer un intervalle de confiance pour  $\theta$  de niveau 95%.