



## CENTRE SCOLAIRE OZANAM

Internat et externat pour lycéens et étudiants – Etudes encadrées et soutien scolaire – Stages intensifs de révision

60 rue Vauban 69006 LYON ☎ 04 78 52 27 99 / Fax : 04 78 52 11 15

✉ [contact@ozanam-lyon.fr](mailto:contact@ozanam-lyon.fr) 🌐 [www.ozanamlyon.fr](http://www.ozanamlyon.fr)

### Concours BCE : Epreuve de Mathématiques II (option Scientifique)

Conception HEC Paris-ESCP Europe : 2 mai 2019

*La régression logique permet de modéliser l'influence qu'exercent des facteurs exogènes sur une variable binaire, c'est-à-dire une variable ne pouvant prendre que deux valeurs.*

*Outre son domaine d'application privilégié qui est l'apprentissage automatique (machine learning), la régression logistique est couramment utilisée aussi bien en médecine qu'en actuariat et en économétrie.*

#### Partie 1. Fonction logistique et lois logistiques

On appelle fonction logistique la fonction  $\Lambda$  définie sur  $\mathbb{R}$  par :  $\forall x \in \mathbb{R}, \Lambda(x) = \frac{1}{1 + e^{-x}}$

1. Montrer que  $\Lambda$  est une bijection dans  $\mathbb{R}$  sur  $]0,1[$ , dont la bijection réciproque est la fonction  $L$  définie par :

$$\forall x \in ]0,1[, L(x) = \ln\left(\frac{x}{1-x}\right)$$

2. Calculer la dérivée de la fonction  $\Lambda$ .
3. Justifier l'existence d'un unique réel  $x_0$  tel que :  $\Lambda(x_0) = x_0$ .
4. Établir pour tout  $x \in \mathbb{R}$ , l'inégalité :  $|\Lambda(x) - x| \leq |x - x_0|$ .

Le script Scilab suivant, dont la ligne (1) définit la fonction  $\Lambda$ , permet de calculer une valeur approchée de  $x_0$  par la méthode de dichotomie.

```
(1) def f('y=Lambda(x)', 'y=1/(1+exp(-x))');  
(2) a=0;  
(3) b=1;  
(4) eps= .....;  
(5) while b-a>eps;  
(6)     c=(a+b)/2;  
(7)     if Lambda(c)>c then .....; else b= .....; end;
```

(8) end;

(9)  $x_0 = (a+b) / 2$

5. Compléter la ligne (7) et justifier le choix des valeurs affectées en lignes (2) et (3) aux variables  $a$  et  $b$ .
6. Quelle valeur maximale peut-on affecter en ligne (4) à la variable  $\text{eps}$  pour être assuré que l'erreur d'approximation commise ne dépasse pas  $10^{-4}$  ?
7. Que peut-on dire de la valeur numérique obtenue par l'instruction (10) suivante ?

(10) `Lambda(x0) - x0`

On note  $\lambda$  la dérivée de la fonction  $\Lambda$ .

8. Vérifier que  $\lambda$  est une densité de probabilité.
9. Préciser la parité de la fonction  $\lambda$  ; donner l'allure de sa courbe représentative dans le plan rapporté à un repère orthogonal et en déterminer les points d'inflexion.

On dit qu'une variable aléatoire  $Z$  suit la *loi logistique standard* si elle admet la fonction  $\lambda$  pour densité.

Pour tout couple  $(r, s) \in \mathbb{R} \times \mathbb{R}_+^*$ , on dit qu'une variable aléatoire  $Y$  suit la loi logistique  $\mathcal{L}(r, s)$  si la variable aléatoire  $Z$  définie par  $Z = \frac{Y - r}{s}$  suit la loi logistique standard.

10. Justifier qu'une variable aléatoire qui suit une loi logistique  $\mathcal{L}(r, s)$  admet des moments de n'importe quel ordre et en indiquer l'espérance.
11. En utilisant la méthode d'inversion, écrire le script d'une fonction Scilab, fonction `S=grandlogis(n, p, r, s)`, fournissant pour tout couple  $(n, p)$  d'entiers strictement positifs, une matrice  $S$  à  $n$  lignes et  $p$  colonnes dont les coefficients sont des simulations de variables aléatoires indépendantes suivant la loi logistique  $\mathcal{L}(r, s)$ .
12. Décrire un procédé permettant de calculer une valeur approchée de la variance de la loi logistique standard à l'aide de la fonction `grandlogis`.

Soit  $U_1$  et  $U_2$  deux variables aléatoires indépendantes suivant chacune la loi exponentielle de paramètre 1.

13. Montrer que la variable aléatoire  $Z = \ln\left(\frac{U_1}{U_2}\right)$  suit la loi logistique standard (on pourra utiliser un changement de variable exponentiel, c'est-à-dire de la forme  $t = e^x$ ).
14. En déduire un nouveau script Scilab permettant de simuler une variable aléatoire suivant la loi logistique standard à l'aide de la fonction `grand`.

## Partie 2. Variance de la loi logistique standard

- Pour tout couple  $(a, b) \in \mathbb{R}^2$ , on note  $\text{Im}(z)$  la partie imaginaire  $b$  du nombre complexe  $z = a + ib$ .

- Pour tout polynôme :  $P = \sum_{k=0}^d a_k X^k \in \mathbb{R}[X]$  de degré  $d \in \mathbb{N}$ , les termes non nuls  $a_k X^k$  sont appelés les monômes  $P$  et les  $a_k$  leurs coefficients.
- Dans la factorisation :  $P = a_d \prod_{k=1}^d (X - z_k)$  de  $P$  dans  $\mathbb{C}[X]$  (lorsque  $d \neq 0$ ), la somme  $\sum_{k=1}^d z_k$  est appelée la somme des racines complexes de  $P$ , que les nombres complexes  $z_1, z_2, \dots, z_d$  soient distincts ou non.

Pour tout  $n \in \mathbb{N}$ , on pose :  $P_n = \sum_{k=0}^n (-1)^k \binom{2n+1}{2k+1} (X-1)^{n-k}$

15. Expliciter les polynômes  $P_0$  et  $P_1$ .

16. Pour tout  $n \in \mathbb{N}^*$ , préciser le degré du polynôme  $P_n$  et donner les coefficients de ses deux monômes de plus hauts degrés.

17. Utiliser le résultat précédent pour montrer que pour tout  $n \in \mathbb{N}^*$ , la somme des racines complexes de  $P_n$  est égale à :

$$\frac{2n(n+1)}{3}$$

Soit  $x \in \mathbb{R}$  et  $n \in \mathbb{N}$ .

18. Justifier les égalités suivantes :

$$\sin((2n+1)x) = \operatorname{Im}((\cos(x) + i \sin(x))^{2n+1}) = \sum_{k=0}^n (-1)^k \binom{2n+1}{2k+1} \cos^{2(n-k)}(x) \times \sin^{2k+1}(x)$$

19. En déduire, pour tout  $x \in ]0, \pi[$ , la relation :

$$\frac{\sin((2n+1)x)}{\sin^{2n+1}(x)} = P_n\left(\frac{1}{\sin^2(x)}\right)$$

20. A l'aide du résultat de la question 17, montrer que pour tout  $n \in \mathbb{N}^*$ , on a :

$$\sum_{k=1}^n \frac{1}{\sin^2\left(\frac{k\pi}{2n+1}\right)} = \frac{2n(n+1)}{3}$$

Soit  $x \in ]0, \frac{\pi}{2}[$ .

21. Justifier les inégalités suivantes :

$$\sin(x) \leq x \leq \tan(x) \text{ et } \frac{1}{\sin^2(x)} - 1 \leq \frac{1}{x^2} \leq \frac{1}{\sin^2(x)}$$

22. En utilisant le résultat de la question 20, en déduire, pour tout  $n \in \mathbb{N}^*$ , l'encadrement :

$$\frac{n(2n-1)}{3} \leq \frac{(2n+1)^2}{\pi^2} \sum_{k=1}^n \frac{1}{k^2} \leq \frac{2n(n+1)}{3}$$

23. Établir l'égalité :

$$\sum_{k=1}^{+\infty} \frac{1}{k^2} = \frac{\pi^2}{6}$$

Soit  $Z$  une variable aléatoire suivant la loi logistique standard.

24. À l'aide d'une intégration par parties, justifier que la variance de  $Z$ , notée  $V(Z)$ , vérifie l'égalité :

$$V(Z) = 4 \int_0^{+\infty} \frac{xe^{-x}}{1+e^{-x}} dx$$

25. Établir pour tout  $n \in \mathbb{N}$ , l'égalité :

$$\int_0^{+\infty} \frac{xe^{-x}}{1+e^{-x}} dx = \sum_{k=0}^n (-1)^k \int_0^{+\infty} xe^{-(k+1)x} dx + I_n$$

$$\text{où } I_n = (-1)^{n+1} \int_0^{+\infty} \frac{xe^{-(n+2)x}}{1+e^{-x}} dx$$

26. Montrer que l'intégrale  $I_n$  tend vers 0 lorsque  $n$  tend vers  $+\infty$  et en déduire l'égalité :

$$\int_0^{+\infty} \frac{xe^{-x}}{1+e^{-x}} dx = \sum_{k=0}^{+\infty} \frac{(-1)^k}{(k+1)^2}$$

27. En utilisant la formule établie dans la question 25, déduire de l'égalité précédente que la variance de  $Z$  est égale à  $\frac{\pi^2}{3}$ .

28. Établir la convergence des deux intégrales :

$$\int_0^{+\infty} \ln(x)e^{-x} dx \text{ et } \int_0^{+\infty} (\ln(x))^2 e^{-x} dx$$

On pose  $I = \int_0^{+\infty} \ln(x)e^{-x} dx$  et  $J = \int_0^{+\infty} (\ln(x))^2 e^{-x} dx$ .

29. En utilisant le résultat de la question 13, calculer  $J - I^2$ .

### Partie 3. Estimation à partir de données binaires

Dans cette partie,  $\theta$  est un paramètre réel inconnu et  $F$  désigne la fonction de répartition d'une variable aléatoire à densité dont une densité  $f$  est continue et strictement positive sur  $\mathbb{R}$ .

Soit  $(Y_n)_{n \in \mathbb{N}^*}$  une suite de variables aléatoires indépendantes définies sur un espace probabilisé  $(\Omega, \mathcal{A}, P_\theta)$  suivant chacune la loi de Bernoulli de paramètre  $F(\theta)$ .

30. Justifier que  $F$  est une bijection de  $\mathbb{R}$  sur  $]0,1[$ . On note  $F^{-1}$  sa bijection réciproque.

Pour tout  $n \in \mathbb{N}^*$ , on pose :  $\bar{Y}_n = \frac{1}{n} \sum_{j=1}^n Y_j$

31. Montrer que la suite  $(\sqrt{n}(\bar{Y}_n - F(\theta)))_{n \in \mathbb{N}^*}$  converge en loi vers une variable aléatoire suivant une loi normale centrée dont on précisera la variance.

Pour tout  $n \in \mathbb{N}^*$  et tout  $\omega \in \Omega$ , on pose  $T_n(\omega) = \begin{cases} F^{-1}(\bar{Y}_n(\omega)) & \text{si } 0 \leq \bar{Y}_n(\omega) < 1 \\ 0 & \text{sinon} \end{cases}$ .

De plus, pour tout  $n \in \mathbb{N}^*$ , on note  $E_n$  l'événement  $[0 < \bar{Y}_n < 1]$ .

32. Calculer  $P_\theta(E_n)$  et trouver la limite de cette probabilité lorsque  $n$  tend vers  $+\infty$ .

Soit  $x \in \mathbb{R}$  et  $n \in \mathbb{N}^*$ .

33. Établir l'égalité ensembliste :

$\{\omega \in E_n / T_n(\omega) \leq x\} = [\bar{Y}_n \leq F(x)] \cap E_n$  et montrer que  $[T_n \leq x]$  est un élément de la tribu  $A$ .

34. Justifier l'encadrement

$$P_\theta([\bar{Y}_n \leq F(x)] \cap E_n) \leq P_\theta([T_n \leq x]) \leq P_\theta([\bar{Y}_n \leq F(x)] \cap E_n) + 1 - P_\theta(E_n)$$

35. Montrer que pour tout  $x \neq 0$ , on a  $\lim_{n \rightarrow +\infty} P_\theta([T_n \leq x]) = \begin{cases} 0 & \text{si } x < \theta \\ 1 & \text{si } x > \theta \end{cases}$ .

36. En déduire que  $(T_n)_{n \in \mathbb{N}^*}$  est une suite convergente d'estimateurs du paramètre  $\theta$ .

Pour tout  $n \in \mathbb{N}^*$  et tout  $\omega \in \Omega$ , on pose :

$$U_n(\omega) = \begin{cases} \frac{T_n(\omega) - \theta}{\bar{Y}_n(\omega) - F(\theta)} & \text{si } \bar{Y}_n(\omega) \neq F(\theta) \\ \frac{1}{f(\theta)} & \text{si } \bar{Y}_n(\omega) = F(\theta) \end{cases}$$

On admet sans démonstration que pour tout  $n \in \mathbb{N}^*$ ,  $U_n$  est une variable aléatoire sur  $(\Omega, A, P_\theta)$ .

Soit  $\varepsilon > 0$ .

Pour tout  $n \in \mathbb{N}^*$ , on note  $B_n(\varepsilon)$  l'événement :  $\left[ \left| U_n - \frac{1}{f(\theta)} \right| \leq \varepsilon \right]$

37. Établir l'existence d'un réel  $\alpha > 0$  tel que :

$$\forall x \in [\theta - \alpha, \theta + \alpha], \left| \frac{1}{f(x)} - \frac{1}{f(\theta)} \right| \leq \varepsilon$$

38. Pour un tel  $\alpha$ , justifier l'inclusion :  $[|T_n - \theta| \leq \alpha] \cap E_n \subset B_n(\varepsilon)$  où  $E_n$  a été défini dans les questions 32 à 36.

39. Montrer que la suite  $(U_n)_{n \in \mathbb{N}^*}$  converge en probabilité vers  $\frac{1}{f(\theta)}$ .

40. En déduire que la suite  $(\sqrt{n}(T_n - \theta))_{n \in \mathbb{N}^*}$  converge en loi vers une variable aléatoire suivant une loi normale centrée dont on précisera la variance.

## Partie 4. Régression logistique

- Dans toute cette partie,  $p$  désigne un entier supérieur ou égal à 2.

- Pour tout couple  $(n, m) \in (\mathbb{N}^*)^2$ , on note  $\mathcal{M}_{n,m}(\mathbb{R})$  l'ensemble des matrices à  $n$  lignes et  $m$  colonnes à coefficients réels et  ${}^tM$  la transposée de toute matrice  $M \in \mathcal{M}_{n,m}(\mathbb{R})$ .
- Pour tout  $m \in \mathbb{N}^*$ , le produit scalaire usuel de deux vecteurs  $u$  et  $v$  de  $\mathbb{R}^m$  est noté  $\langle u, v \rangle$ . Si  $U$  et  $V$  sont les matrices colonnes représentant  $u$  et  $v$  dans la base canonique, le produit scalaire  $\langle u, v \rangle$  est donc l'unique coefficient de la matrice  ${}^tUV$ .
- On rappelle que les fonctions  $\Lambda$  et  $L$  ont été définies dans la partie 1.

Dans cette partie, on note  $Y$  une variable aléatoire de Bernoulli, dite *variable endogène*, dont la loi dépend du niveau de  $p$  *facteurs exogènes*.

L'influence de ces facteurs sur la loi de  $Y$  est résumée par la fonction  $b$  qui associe à un vecteur  $x \in \mathbb{R}^p$ , la probabilité  $b(x)$  que  $Y$  soit égale à 1 lorsque les niveaux des facteurs sont donnés par les composantes du vecteur  $x$ .

Dans le *modèle de régression logistique* envisagé dans cette partie, la fonction  $b$  est supposée de la forme :

$$b : x \mapsto \Lambda(\langle \alpha, x \rangle)$$

où  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_p)$  est un vecteur de  $\mathbb{R}^p$  dont les composantes  $\alpha_1, \alpha_2, \dots, \alpha_p$  sont des paramètres inconnus qui représentent les degrés d'influence des divers facteurs exogènes sur la variable endogène  $Y$ .

Pour estimer les paramètres du modèle, on dispose de  $k$  vecteurs  $x^{(1)}, x^{(2)}, \dots, x^{(k)}$  de  $\mathbb{R}^p$  ( $k \in \mathbb{N}^*$ ) et pour tout  $i \in \llbracket 1, k \rrbracket$ , d'une suite  $(Y_{i,n})_{n \in \mathbb{N}^*}$  de variables aléatoires indépendantes suivant chacune la loi de Bernoulli de paramètre  $b(x^{(i)}) = \Lambda(\langle \alpha, x^{(i)} \rangle)$ .

Pour chaque indice fixé  $i$  et pour tout  $n \in \mathbb{N}^*$ , les variables aléatoires  $Y_{i,1}, Y_{i,2}, \dots, Y_{i,n}$  définissent donc un  $n$ -échantillon associé à la loi endogène lorsque les niveaux exogènes sont les composantes  $x_1^{(i)}, x_2^{(i)}, \dots, x_p^{(i)}$  du vecteur  $x^{(i)}$  dans la base canonique de  $\mathbb{R}^p$ .

On note respectivement  $A$  et  $M$  la matrice du vecteur  $\alpha$  et la matrice de la famille  $(x^{(1)}, x^{(2)}, \dots, x^{(k)})$  dans la base canonique de  $\mathbb{R}^p$  :

$$A = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_p \end{pmatrix} \in \mathcal{M}_{p,1}(\mathbb{R}) \quad \text{et} \quad M = \begin{pmatrix} x_1^{(1)} & \dots & x_1^{(k)} \\ \vdots & & \vdots \\ x_p^{(1)} & \dots & x_p^{(k)} \end{pmatrix} \in \mathcal{M}_{p,k}(\mathbb{R})$$

On suppose que le rang de la matrice  $M$  est égale à  $p$ .

41. Montrer que la matrice  $M^t M$  est inversible.

42. Montrer que pour toute matrice  $H \in \mathcal{M}_{k,1}(\mathbb{R})$ , la matrice  $U \in \mathcal{M}_{p,1}(\mathbb{R})$  pour laquelle l'unique coefficient de la matrice  ${}^t(tMU - H)({}^tMU - H)$  est le plus petit possible, est la matrice  $(M^t M)^{-1} M H$ .

43. Expliquer pourquoi les lois des variables aléatoires  $Y_{i,n}$  ne suffiraient pas à définir le vecteur  $\alpha$  si le rang de  $M$  n'était pas égal à  $p$ .

Pour tout  $n \in \mathbb{N}^*$  et tout  $i \in \llbracket 1, k \rrbracket$ , on pose :  $\bar{Y}_{i,n} = \frac{1}{n} \sum_{j=1}^n Y_{i,j}$  et pour tout  $\omega \in \Omega$  :

$$T_{i,n} = \begin{cases} L(\bar{Y}_{i,n}(\omega)) & \text{si } 0 < \bar{Y}_{i,n}(\omega) < 1 \\ 0 & \text{sinon} \end{cases}$$

Soit  $(c_1, c_2, \dots, c_k) \in \mathbb{R}^k$ .

44. En utilisant les résultats de la partie 3, montrer que :

$\left( \sum_{i=1}^k c_i T_{i,n} \right)_{n \in \mathbb{N}^*}$  est une suite convergente d'estimateurs du paramètre  $\sum_{i=1}^k c_i \langle \alpha, x^{(i)} \rangle$

Pour tout  $n \in \mathbb{N}^*$  et tout  $\omega \in \Omega$ , on pose :

$$H_n(\omega) = \begin{pmatrix} T_{1,n}(\omega) \\ T_{2,n}(\omega) \\ \vdots \\ T_{k,n}(\omega) \end{pmatrix} \text{ et } \begin{pmatrix} A_{1,n}(\omega) \\ A_{2,n}(\omega) \\ \vdots \\ A_{p,n}(\omega) \end{pmatrix} = (M^t M)^{-1} M H_n(\omega)$$

45. Montrer que pour tout  $j \in \llbracket 1, p \rrbracket$ , la suite  $(A_{j,n})_{n \in \mathbb{N}^*}$  est une suite convergente d'estimateurs de  $\alpha_j$ .

www.ozanamlyon.fr