

Statistique

Estimation sous biais de sélection et avec fonction de poids inconnue

Agathe Guilloux

LSTA, université Pierre et Marie Curie, 175, rue du Chevaleret, 75013 Paris, France

Reçu le 3 février 2005 ; accepté après révision le 15 novembre 2005

Présenté par Paul Deheuvels

Résumé

Nous considérons le problème de l'estimation de la fonction de répartition G d'une variable aléatoire (v.a.) positive X à partir de l'observation d'une v.a. biaisée Y de fonction de répartition $F_w = \int w(x) dG(x) / \mu_w$, où w est une fonction de poids inconnue. En supposant de plus que l'échantillon issu de la fonction de répartition F_w est censuré à droite, nous construisons un estimateur \widehat{G} de la fonction de répartition G pour lequel on énonce un théorème de consistance forte et de convergence faible. **Pour citer cet article :** A. Guilloux, C. R. Acad. Sci. Paris, Ser. I 342 (2006).

© 2005 Académie des sciences. Publié par Elsevier SAS. Tous droits réservés.

Abstract

Distribution estimation from biased data with unknown weighting function. We consider the problem of estimating the cumulative distribution function (cdf) G of a non-negative random variable (r.v.) X from the observation of a biased r.v. Y with cdf $F_w = \int w(x) dG(x) / \mu_w$, where w is an unknown weighting function. We assume moreover that the random sample with common cdf F_w is right-censored. We construct an estimator \widehat{G} for the cdf G and state its strong consistency and weak convergence. **To cite this article :** A. Guilloux, C. R. Acad. Sci. Paris, Ser. I 342 (2006).

© 2005 Académie des sciences. Publié par Elsevier SAS. Tous droits réservés.

1. Introduction

Considérons une population d'individus $i \in I$, dans laquelle on note σ_i l'instant de naissance et X_i la durée de vie de l'individu i . Soit G la fonction de répartition de la v.a. X . Supposons que l'on ne peut observer que les individus vivants à un instant t_0 fixé d'échantillonnage. L'individu i peut donc entrer dans l'échantillon si $\sigma_i < t_0$ et $\sigma_i + X_i > t_0$ (ce qui assure qu'il est né avant t_0 et qu'il meurt après t_0). On peut alors montrer que les v.a. ζ et Y , qui représentent, respectivement, l'instant de naissance et la durée de vie pour les individus vivants à l'instant t_0 , souffrent d'un biais de sélection. On a, en effet, pour tout s et t :

$$\mathbb{P}(\zeta \leq s, Y \leq t) = \mathbb{P}(\sigma \leq s, X \leq t | \sigma < t_0, \sigma + X > t_0) \neq \mathbb{P}(\sigma \leq s, X \leq t).$$

Adresse e-mail : aguillou@ccr.jussieu.fr (A. Guilloux).

Plus précisément, supposons que le processus ponctuel $\eta = \sum_{i \in I} \delta_{\sigma_i}$ (où δ_a est la masse de Dirac en a), formé à partir des instants de naissance dans la population I , est poissonnien non-homogène d'intensité φ . On peut alors montrer, voir en particulier Lund [8], que Y a pour fonction de répartition F_w donnée, pour tout $t > 0$, par :

$$F_w(t) = \frac{\int_0^t w(x) dG(x)}{\int_0^\infty w(x) dG(x)}, \quad \text{où } w(t) = \int_0^t \varphi(t_0 - \sigma) d\sigma, \quad \text{pour tout } t > 0, \quad (1)$$

et $\mu_w = \int_0^\infty w(x) dG(x) < \infty$. La fonction F_w est souvent appelée version *biaisée* de la fonction G et l'on dit alors que la v.a. X souffre d'un biais de sélection.

Depuis Fisher [5], de nombreux auteurs ont étudié ce problème, en particulier Patil et Rao [9], Gill et al. [6] et Efromovitch [4]. Ces derniers ont fait l'hypothèse que la fonction de poids w est connue. Notre but est ici de montrer que, dans le modèle d'échantillonnage décrit ci-dessus, on peut se passer de cette hypothèse. Dans ce cadre et dans le cas où la fonction de biais w est inconnue, nous construisons un estimateur \widehat{G} de la fonction de répartition G à partir d'un n -échantillon de couples $(\zeta_1, Y_1), \dots, (\zeta_n, Y_n)$, où les durées de vie Y_1, \dots, Y_n sont censurées à droite par des v.a. positives C_1, \dots, C_n . Nous énonçons également un théorème de consistance uniforme pour l'estimateur \widehat{G} et de convergence faible pour le processus $\sqrt{n}(\widehat{G} - G)$.

2. Description des observations et de la censure

Considérons maintenant la population J des individus vivants à l'instant t_0 . Soit un individu j de cette population J , né en ζ_j et de durée de vie Y_j . La v.a. positive $t_0 - \zeta_j$ représente alors son âge à l'instant t_0 d'échantillonnage et peut donc être appelée « temps de récurrence arrière ». On peut montrer, toujours en supposant que le processus $\eta = \sum_{i \in I} \delta_{\sigma_i}$ est poissonnien non-homogène d'intensité φ , que, pour tout $ut \geq 0$:

$$\mathbb{P}(t_0 - \zeta \leq t) = \mathbb{P}(t_0 - \sigma \leq t | \sigma < t_0, \sigma + X > t_0) = \frac{1}{\mu_w} \int_0^{t \wedge t_0} \varphi(t_0 - s) \overline{G}(s) ds. \quad (2)$$

Ainsi la loi des v.a. ζ et Y , voir les Éqs. (1) et (2), détermine de manière unique les fonctions G et w . Ceci assure l'identifiabilité du problème posé.

La durée de vie après l'échantillonnage de l'individu j est donnée par la v.a. positive $\zeta_j + y_j - t_0$ qui peut alors être nommée « temps de récurrence avant ». Comme les individus entrent dans l'échantillon au temps t_0 , il est naturel de supposer que seul le temps de récurrence avant peut être censuré ou, de façon équivalente, qu'un individu ne peut être censuré qu'à partir du moment où il est dans l'échantillon. On peut alors supposer que les durées de vie des individus de la population J sont censurées à droite par une v.a. C de fonction de répartition H de telle sorte que l'on n'observe pas des réalisations de la v.a. Y mais de la v.a. Z donnée par :

$$Z = t_0 - \zeta + (\zeta + Y - t_0) \wedge C$$

où $x \wedge y$ est le minimum entre x et y . On suppose également que l'indicatrice $I(\{\zeta + Y - t_0 \leq C\})$ est observable, où $I(A)$ est l'indicatrice de l'événement A .

Ce modèle de censure a été étudié par Winter et Földes [11] et Asgharian et al. [3], notamment, et défendu vivement par Asgharian [2]. Il faut noter ici que ces derniers auteurs ont étudié le cas particulier où, dans l'expression de la fonction de répartition F_w , la fonction de poids est donnée par $w(t) = t$ pour tout $t > 0$, ce cas particulier est appelé *biais de longueur*. Ceci correspond, dans notre formulation, à supposer que le processus φ est poissonnien homogène, c'est-à-dire d'intensité φ constante.

3. Estimation de la fonction de répartition

Comme décrit dans la partie précédente, nous travaillons maintenant avec n individus vivants à l'instant t_0 pour lesquels on observe les triplets indépendants :

$$(\zeta_j, Z_j, I(\{\zeta_j + Y_j - t_0 \leq C_j\})) \quad \text{pour } j = 1, \dots, n.$$

Download English Version:

<https://daneshyari.com/en/article/4671880>

Download Persian Version:

<https://daneshyari.com/article/4671880>

[Daneshyari.com](https://daneshyari.com)