

Introduction au bootstrap

Résumé

Présentation succincte du principe du bootstrap.

Retour au [plan du cours](#)

1 Introduction

La motivation du *bootstrap*¹ (Efron, 1982 ; Efron et Tibshirani, 1993) est d’approcher par simulation (*Monte Carlo*) la distribution d’un estimateur lorsque l’on ne connaît pas la loi de l’échantillon ou, plus souvent lorsque l’on ne peut pas supposer qu’elle est gaussienne. L’objectif est de remplacer des hypothèses probabilistes pas toujours vérifiées ou même invérifiables par des simulations et donc beaucoup de calcul.

Le principe fondamental de cette technique de ré-échantillonnage est de substituer à la distribution de probabilité inconnue F , dont est issu l’échantillon d’apprentissage, la distribution empirique \hat{F} qui donne un poids $1/n$ à chaque réalisation. Ainsi on obtient un échantillon de taille n dit *échantillon bootstrap* selon la distribution empirique \hat{F} par n tirages aléatoires avec remise parmi les n observations initiales.

Il est facile de construire un grand nombre d’échantillons bootstrap sur lesquels calculer l’estimateur concerné. La loi simulée de cet estimateur est une approximation asymptotiquement convergente sous des hypothèses raisonnables² de la loi de l’estimateur. Cette approximation fournit ainsi des estimations du biais, de la variance, donc d’un risque quadratique, et même des intervalles de confiance de l’estimateur sans hypothèse (normalité) sur la vraie loi.

1. Cette appellation est inspirée du baron de Münchhausen (Rudolph Erich Raspe) qui se sortit de sables mouvants par traction sur ses *tirants de bottes*. En France “bootstrap” est parfois traduit par *à la Cyrano* (acte III, scène 13) en référence à ce héros qui prévoyait d’atteindre la lune en se plaçant sur une plaque de fer et en itérant le jet d’un aimant.

2. Échantillon indépendant de même loi et estimateur indépendant de l’ordre des observations.

1.1 Principe du *plug-in*

Soit $\mathbf{x} = \{x_1, \dots, x_n\}$ un échantillon de taille n issue d’une loi inconnue F sur (Ω, \mathcal{A}) . On appelle *loi empirique* \hat{F} la loi discrète des singletons (x_1, \dots, x_n) affectés des poids $1/n$:

$$\hat{F} = \sum_{i=1}^n \delta_{x_i}.$$

Soit $A \in \mathcal{A}$, $P_F(A)$ est estimée par :

$$(\hat{P})_F(A) = P_{\hat{F}}(A) = \sum_{i=1}^n \delta_{x_i}(A) = \frac{1}{n} \text{Card} x_i \in A.$$

De manière plus générale, soit θ un paramètre dont on suppose que c’est une fonction de la loi F . on écrit donc $\theta = t(F)$. Par exemple, $\mu = E(F)$ est un paramètre de F suivant ce modèle. Une *statistique* est une fonction (mesurable) de l’échantillon. Avec le même exemple :

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

et \bar{x} est la statistique qui estime μ . On dit que c’est un estimateur “plug-in” et, plus généralement,

DÉFINITION 1. — On appelle *estimateur plug-in* d’un paramètre θ de F , l’estimateur obtenu en remplaçant la loi F par la loi empirique :

$$\hat{\theta} = t(\hat{F}).$$

comme dans le cas de l’estimation de μ : $\hat{\mu} = E(\hat{F}) = \bar{x}$.

1.2 Estimation de l’écart-type de la moyenne

Soit X une variable aléatoire réelle de loi F . On pose :

$$\mu_F = E_F(X), \quad \text{et} \quad \sigma_F^2 = \text{Var}_F(X) = E_F[(X - \mu_F)^2];$$

Ce qui s'écrit :

$$X \sim (\mu_F, \sigma_F^2).$$

Soit (X_1, \dots, X_n) n variables aléatoires i.i.d. suivant aussi la loi F . Posons $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$. Cette variable aléatoire a pour espérance μ_F et pour variance σ_F^2/n . On dit aussi que la statistique

$$\bar{X} \sim (\mu_F, \sigma_F^2/n).$$

Remarquons qu'en moyennant plusieurs valeurs ou observations, on réduit la variance inhérente à une observation. De plus, sous certaines conditions sur la loi F et comme résultat du théorème de la limite centrale, \bar{X} converge en loi vers la loi normale.

L'estimateur plug-in de σ_F est défini par :

$$\begin{aligned} \hat{\sigma}^2 &= \hat{\sigma}_F^2 = \sigma_{\hat{F}}^2 = \text{Var}_{\hat{F}}(X) \\ &= E_{\hat{F}}[(X - E_{\hat{F}}(X))^2] = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2. \end{aligned}$$

L'estimateur plug-in de σ_F est (légèrement) différent de celui du maximum de vraisemblance. L'estimateur plug-in est en général biaisé mais il a l'avantage d'être simple et de pouvoir s'appliquer à tout paramètre θ même lorsque l'on ne peut pas calculer la vraisemblance du modèle.

2 Estimation bootstrap d'un écart-type

Soit $\hat{\theta} = s(x)$ un estimateur quelconque (M.V. ou autre) de θ pour un échantillon x donné. On cherche à apprécier la précision de $\hat{\theta}$ et donc à estimer son écart-type.

2.1 Échantillon bootstrap

Avec les mêmes notations, \hat{F} est la distribution empirique d'un échantillon $\mathbf{x} = \{x_1, \dots, x_n\}$.

DÉFINITION 2. — On appelle échantillon bootstrap de \mathbf{x} un échantillon de taille n noté

$$\mathbf{x}^* = \{x_1^*, \dots, x_n^*\}$$

suivant la loi \hat{F} ; \mathbf{x}^* est un ré-échantillon de \mathbf{x} avec remise.

2.2 Estimation d'un écart-type

DÉFINITION 3. — On appelle estimation bootstrap de l'écart-type $\hat{\sigma}_F(\hat{\theta})$ de $\hat{\theta}$, son estimation plug-in : $\sigma_{\hat{F}}(\hat{\theta})$.

Mais, à part dans le cas très élémentaire où, comme dans l'exemple ci-dessus, θ est une moyenne, il n'y a pas de formule explicite de cet estimateur. Une approximation de l'estimateur bootstrap (ou plug-in) de l'écart-type de $\hat{\theta}$ est obtenue par une simulation (Monte-Carlo) décrite dans l'algorithme ci-dessous.

Pour un paramètre θ et un échantillon \mathbf{x} donnés, on note $\hat{\theta} = s(\mathbf{x})$ l'estimation obtenue sur cet échantillon. Une *réplication bootstrap* de $\hat{\theta}$ est donnée par : $\hat{\theta}^* = s(\mathbf{x}^*)$.

ALGORITHME 1 : Estimation de l'écart-type

Soit \mathbf{x} un échantillon et θ un paramètre.

for $b = 1$ à B **do**

Sélectionner 1 échantillon bootstrap $\mathbf{x}^{*b} = \{x_1^{*b}, \dots, x_n^{*b}\}$. par tirage avec remise dans \mathbf{x} .

Estimer sur cet échantillon : $\hat{\theta}^*(b) = s(\mathbf{x}^{*b})$.

end for

Calculer l'écart-type de l'échantillon ainsi construit :

$$\begin{aligned} \hat{\sigma}_B^2 &= \frac{1}{B-1} \sum_{b=1}^B (\hat{\theta}^*(b) - \hat{\theta}^*(.))^2 \\ \text{avec } \hat{\theta}^*(.) &= \frac{1}{B} \sum_{b=1}^B \hat{\theta}^*(b). \end{aligned}$$

$\hat{\sigma}_B$ est l'approximation bootstrap de l'estimation plug-in recherchée de l'écart-type de $\hat{\theta}$.

2.3 Estimation du biais

Avec les mêmes notations :

$$\theta = t(F) \quad \text{et} \quad \hat{\theta} = s(\mathbf{x}),$$

le biais d'un estimateur s'exprime comme

$$\mathcal{B}_F(\hat{\theta}) = E_F[s(\mathbf{x})] - t(F).$$

Un estimateur est sans biais si $E[\hat{\theta}] = \theta$. Le biais est aussi une mesure de la précision d'un estimateur et on a vu que, généralement, les estimateurs plug-in étaient biaisés.

DÉFINITION 4. — *On appelle estimateur bootstrap du biais, l'estimateur plug-in :*

$$\widehat{\mathcal{B}}_F(\hat{\theta}) = \mathcal{B}_{\widehat{F}}(\hat{\theta}) = E_{\widehat{F}}[s(\mathbf{x}^*)] - t(\widehat{F}).$$

Comme pour l'écart-type, il n'existe généralement pas d'expression analytique et il faut avoir recours à une approximation par simulation.

ALGORITHME 2 : *Estimation bootstrap du biais*

Soit \mathbf{x} un échantillon et θ un paramètre.

for $b = 1$ à B **do**

*Sélectionner 1 échantillon bootstrap $\mathbf{x}^{*b} = \{x_1^{*b}, \dots, x_n^{*b}\}$. par tirage avec remise dans \mathbf{x} .*

Estimer sur cet échantillon la réplique bootstrap de $\hat{\theta}$: $\hat{\theta}^(b) = s(\mathbf{x}^{*b})$.*

end for

Approcher $E_{\widehat{F}}[s(\mathbf{x}^)]$ par $\hat{\theta}^*(.) = \frac{1}{B} \sum_{b=1}^B (\hat{\theta}^*(b))$*

L'approximation bootstrap du biais est : $\widehat{\mathcal{B}}_B(\hat{\theta}) = \hat{\theta}^(.) - \hat{\theta}$.*

Le bootstrap rapidement décrit ici est dit “non-paramétrique” car la loi empirique \widehat{F} est une estimation non-paramétrique de F . Dans le cas où F serait connue à un paramètre près, il existe également une version dite *paramétrique* du bootstrap.

Pour des estimateurs plus compliqués (fonctionnels) comme dans le cas de la régression non-paramétrique par noyau ou spline, il est facile de construire graphiquement une enveloppe bootstrap de l'estimateur à partir de répliques de l'échantillon. Celle-ci fournit généralement une bonne appréciation de la qualité de l'estimateur obtenu. Attention, dans le cas de la régression il est en principe plus justifié de répliquer le tirage sur les *résidus* plutôt que sur les observations. Ce sont les résidus qui sont en effet supposés i.i.d. et qui vérifient donc les hypothèses nécessaires mais cette approche devient très sensible à l'hypothèse sur la validité du modèle. Il est finalement d'usage de considérer un échantillon bootstrap issu des données initiales (Efron et Tibshirani) :

$$\mathbf{z}^{*b} = \{(\mathbf{x}_1^{*b}, y_1^{*b}), \dots, (\mathbf{x}_n^{*b}, y_n^{*b})\};$$

c'est ce qui a été choisi dans ce document.

Enfin, l'estimation bootstrap est justifiée par des propriétés asymptotiques (convergence en loi) lorsque le nombre de répliques (B) croît conjointement avec la taille de l'échantillon (n). Comme la loi empirique \widehat{F} converge (en loi) vers celle théorique, la distribution du paramètre $\hat{\theta} = t(\widehat{F})$ converge (en loi) vers celle théorique de $\theta = t(F)$.

3 Compléments

En résumé, on peut dire que le bootstrap repose sur une hypothèse très élémentaire : $\hat{\theta}^*$ se comporte par rapport à $\hat{\theta}$ comme $\hat{\theta}$ par rapport à θ . La connaissance de $\hat{\theta}^*$ (distribution, variance, biais...) renseigne alors sur celle de $\hat{\theta}$.

Beaucoup d'autres compléments sont à rechercher dans la littérature et en particulier dans Efron et Tibshirani (1993). Il est ainsi possible de définir des intervalles de confiance bootstrap en considérant la distribution et les quantiles de $\hat{\theta}^*$ ou même encore des tests à partir des versions bootstrap de leur statistique.