

Méthodes de ré-échantillonnage bootstrap

Master 1 Mathématiques Appliquées
Université Paris-Saclay

28 mars 2019

Qu'est-ce que le bootstrap ?

"To pull oneself up by one's own bootstraps"

Le bootstrap est une méthode statistique pour évaluer la précision d'une méthode : estimer une variance, une erreur de prédiction, calculer un intervalle de confiance,...

Pour estimer la variabilité d'un estimateur (sa loi), il faut observer plusieurs réalisations.

↔ une seule expérience en général

↔ Idée : générer des échantillons qui ressemblent à l'échantillon de départ !

$T(Y_1, \dots, Y_n)$ une statistique basée sur n observations i.i.d.
Quelle est sa loi ?

↪ calcul (asymptotique), souvent difficile

↪ estimation par simulation

La loi des Y_i n'est pas connue, il faut l'estimer : $\hat{P} = P_{\hat{\theta}(Y)}$.

↪ hyp. *robustesse* et *régularité* du modèle.

Contexte **non-paramétrique** :

$$P_n = \frac{1}{n} \sum_i \delta_{Y_i}$$

P_n est la mesure empirique (la loi uniforme sur l'échantillon), c'est un bon estimateur de P (th. de Glivenko-Cantelli).

Méthode empirique : soit $\theta = T(P)$ un paramètre ; on estime $T(P)$ par $T(P_n)$.

Exemples : $\mu(P) = \int x dP(x)$ ou $\sigma^2(P) = \int x^2 dP(x) - (\mu(P))^2$.

La donnée de l'échantillon (Y_1, \dots, Y_n) est équivalente à celle de P_n . On interprète la statistique comme une fonction de P_n :

$$\hat{\theta} = T(P_n),$$

est un estimateur de $T(P)$.

Exemple : $\sigma^2(P) = E_P[(Y - \mu(P))^2]$ a pour estimateur

$$\sigma^2(P_n) = E_{P_n}[(Y - \mu(P_n))^2] = \frac{1}{n} \sum_i (Y_i - \bar{Y})^2$$

C'est un échantillon (en général de taille n) tiré selon P_n parmi Y_1, \dots, Y_n .
On le note $Y^* = (Y_1^*, \dots, Y_n^*)$.

- Dans le monde réel, on estime $\theta = T(P)$ par $\hat{\theta} = T(P_n) = T(Y_1, \dots, Y_n)$;
- Dans le monde bootstrap, on estime θ par $T(Y_1^*, \dots, Y_n^*)$.

Le bootstrap est valide si on peut montrer que la loi de $T(Y^*) = \hat{\theta}^*$ est une bonne approximation de la loi de $T(Y) = \hat{\theta}$, c'est à dire si $\mathcal{L}_P(\hat{\theta})$ et $\mathcal{L}_{P_n}(\hat{\theta}^* | Y_1, \dots, Y_n)$ ont même loi limite.

Conditions générales de validité :

- la loi limite de $T(P_n)$ dépend "continûment" de P
- la convergence est uniforme en P dans un voisinage de la vraie loi P_0
- la loi utilisée pour générer l'échantillon bootstrap estime bien P_0

Théorème (Bickel et Freedman, Annals of Statistics 1981) :

si $\theta = \int YdP$,

$$\lim_{n \rightarrow \infty} \sqrt{n} \sup_x |P(T_n \leq x) - P(T_n^* \leq x | Y_1, \dots, Y_n)| = 0$$

p.s. conditionnellement à (Y_1, \dots, Y_n) , avec $T_n = \sqrt{n}(\hat{\theta} - \theta)$ et $T_n^* = \sqrt{n}(\hat{\theta}^* - \hat{\theta})$.

Ce résultat a été étendu aux fonctions de la moyenne, modèles de régression, séries temporelles...

Y_1, \dots, Y_n de loi uniforme $(0, \theta)$.

$$\hat{\theta} = \max(Y_1, \dots, Y_n)$$

$\hookrightarrow n(\hat{\theta} - \theta)/\theta$ de loi exponentielle $\mathcal{E}(1)$

$\hookrightarrow \Pr(\hat{\theta}^* = \hat{\theta}) \simeq 0.632$. *Pourquoi?*

En pratique : une méthode bootstrap doit toujours être validée par un théorème, ou par une étude de simulations.

Ré-échantillonnage bootstrap

$Y = (Y_1, \dots, Y_n)$ échantillon initial

Échantillon bootstrap : n observations tirées uniformément avec remplacement parmi les observations de départ

$$\hookrightarrow Y^* = Y_1^*, Y_2^*, \dots, Y_n^*$$

Problème : en général, on ne connaît pas mieux la loi de $T(Y^*)$ que celle de $T(Y)$!

Avantage : on peut générer autant d'échantillons Y^* que l'on veut!

$$\hookrightarrow Y^{*1} = (Y_1^{*1}, \dots, Y_n^{*1})$$

⋮

$$\hookrightarrow Y^{*B} = (Y_1^{*B}, \dots, Y_n^{*B})$$

En théorie, $B = n^n$; en pratique $B = 200$ à 1000 .

C'est l'étape de Monte-Carlo :

$$\begin{array}{ccc} (Y_1^{*(1)}, \dots, Y_n^{*(1)}) & \rightarrow & \hat{\theta}_n^{*(1)} \\ & \dots & \dots \\ (Y_1^{*(b)}, \dots, Y_n^{*(b)}) & \rightarrow & \hat{\theta}_n^{*(b)} \\ & \dots & \dots \\ (Y_1^{*(B)}, \dots, Y_n^{*(B)}) & \rightarrow & \hat{\theta}_n^{*(B)} \end{array}$$

Il y a donc deux étapes d'approximation dans le bootstrap :

- 1 l'approximation statistique (asymptotique) ($n \rightarrow \infty$)
- 2 l'approximation numérique (Monte-Carlo) (B grand).

$$\hat{\theta}_n^{*(1)}, \hat{\theta}_n^{*(2)}, \dots, \hat{\theta}_n^{*(B)}$$

est appelée **loi d'échantillonnage bootstrap** de $T(Y) = \hat{\theta}$.

Dans un cadre paramétrique, si $\hat{\theta} = T(\mathbf{Y})$ estime θ :

$$\begin{aligned} \text{Loi } (\hat{\theta} - \theta) &\approx \text{Loi bootstrap}(\hat{\theta}^* - \hat{\theta}) \\ &\approx \text{Loi d'échantillonnage } (\hat{\theta}^{*1} - \hat{\theta}, \hat{\theta}^{*2} - \hat{\theta}, \dots, \hat{\theta}^{*B} - \hat{\theta}) \end{aligned}$$

- Dans certains cas, l'approximation bootstrap est meilleure que l'approximation asymptotique gaussienne de la loi de l'estimateur du Maximum de Vraisemblance

- Permet d'apprécier par simulation la variabilité donc la précision d'un estimateur
- Simplicité du principe et de la mise en oeuvre : pas de formules mathématiques compliquées
- Méthode très générale dans un cadre paramétrique ou non-paramétrique (sous réserve de validité)
- le bootstrap "de base" est valide typiquement dans des situations où l'approximation gaussienne est valide.
Contre exemples : valeurs extrêmes (max, quantiles), U-statistiques. Il existe alors des solutions par sous-échantillonnage.
- le bootstrap présenté est le bootstrap *non-paramétrique*. Il existe des bootstraps paramétriques.

Le biais de $\hat{\theta} = \theta(P_n)$ est :

$$B_P(\hat{\theta}, \theta) = E_P(\hat{\theta}) - \theta.$$

La contrepartie bootstrap du biais est

$$B_{P_n}(\hat{\theta}^*, \hat{\theta}) = E_{P_n}(\hat{\theta}^*) - \hat{\theta}.$$

Le biais est estimé par bootstrap par

$$\widehat{\text{Biais}}^* = \frac{1}{B} \sum_{b=1}^{b=B} \hat{\theta}^{*(b)} - \hat{\theta}.$$

- 1 Générer $Y_1^*, \dots, Y_n^* \sim P_n$
- 2 Calculer $T^* = T(Y_1^*, \dots, Y_n^*)$
- 3 Répéter les étapes 1 et 2 B fois pour obtenir $T^{*(1)}, T^{*(2)}, \dots, T^{*(B)}$
- 4 Estimateur de la variance de la loi de T_n :

$$V_{\text{boot}} = \frac{1}{B} \sum_{b=1}^{b=B} \left(T^{*(b)} - \frac{1}{B} \sum_{b=1}^{b=B} T^{*(b)} \right)^2$$

Principe : on approche la loi de $\hat{\theta} - \theta$ par la loi de $\hat{\theta}^* - \hat{\theta}$.

↔ on calcule les quantiles bootstrap $b_{\alpha/2}^*$ et $b_{1-\alpha/2}^*$ de la loi d'échantillonnage bootstrap

$$(\hat{\theta}^{*1} - \hat{\theta}, \hat{\theta}^{*2} - \hat{\theta}, \dots, \hat{\theta}^{*B} - \hat{\theta})$$

L'intervalle bootstrap de niveau (asymptotique) $1 - \alpha$ est

$$IC(\theta) = [\hat{\theta} - b_{1-\alpha/2}^*; \hat{\theta} - b_{\alpha/2}^*]$$

Pourquoi?

$\{(x_i, y_i), i = 1, \dots, n\}$ données d'apprentissage

$$y_i = g(x_i) + \varepsilon_i$$

- ↪ bootstrap des résidus $\hat{\varepsilon}_i = Y_i - \hat{g}(x_i)$
- ↪ bootstrap des couples
- ↪ bootstrap paramétrique des résidus

- *An Introduction to Statistical Learning, with Applications in R*, Gareth, Witten, Hastie and Tibshirani
- *Bootstrap Methods and their Application*, Davison and Hinkley