

Bootstrap

Master parcours SSD - UE Statistique Computationnelle

Pierre Mahé - bioMérieux & Université de Grenoble-Alpes

- ▶ **Bootstrap** : méthode d'inférence basée sur le **ré-échantillonnage** d'un échantillon.
- ▶ Approche **non-paramétrique** : pas d'hypothèse sur la loi de la variable aléatoire sous-jacente.
- ▶ Principe générique décliné pour **différentes applications**.

- ▶ Introduction
- ▶ Formalisation du principe de **ré-échantillonnage**
- ▶ Le bootstrap pour l'**inférence statistique** :
 - ▶ caractérisation d'un estimateur
 - ▶ intervalles de confiance
- ▶ **2 applications** classiques :
 - ▶ régression
 - ▶ construction de modèles de prédiction
- ▶ Conclusion

Introduction

Formalisation

Bootstrap pour
l'inférence

Caractérisation
d'un estimateur
Intervalle de
confiance

Applications

Bootstrap &
régression

Bootstrap &
prédiction

Conclusion

Références

Introduction

Cours précédent :

- ▶ méthodes de **simulation** pour répondre à des questions d'**inférence statistique**.
- ▶ données simulées selon des **modèles paramétriques**.

Cours précédent :

- ▶ méthodes de **simulation** pour répondre à des questions d'**inférence statistique**.
- ▶ données simulées selon des **modèles paramétriques**.

Le bootstrap :

- ▶ même objectif mais vise à **relâcher ces hypothèses**.
- ▶ se base uniquement sur le vecteur d'observations disponibles : **ré-échantillonnage**.

⇒ une approche totalement **non paramétrique**

Le bootstrap en deux mots

On s'intéresse à un estimateur $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$ d'un paramètre θ .

Outline

UE StatComp

Introduction

Formalisation

Bootstrap pour l'inférence

- Caractérisation d'un estimateur
- Intervalle de confiance

Applications

- Bootstrap & régression
- Bootstrap & prédiction

Conclusion

Références

Le bootstrap en deux mots

On s'intéresse à un estimateur $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$ d'un paramètre θ .

On veut caractériser sa distribution d'échantillonnage, **en se basant uniquement sur l'échantillon (X_1, \dots, X_n)** .

Le bootstrap en deux mots

On s'intéresse à un estimateur $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$ d'un paramètre θ .

On veut caractériser sa distribution d'échantillonnage, **en se basant uniquement sur l'échantillon (X_1, \dots, X_n)** .

On applique la procédure suivante :

- ▶ Pour b allant de 1 à B ,
- ▶ On génère un échantillon (X_1^*, \dots, X_n^*) en **tirant avec remise** dans (X_1, \dots, X_n) .
- ▶ On calcule la valeur de notre estimateur $\hat{\theta}^{*(b)}$ à partir de (X_1^*, \dots, X_n^*) .

Et on travaille sur les B réalisations $(\hat{\theta}^{*(b)})_{b=1, \dots, B}$.

Le bootstrap en deux mots

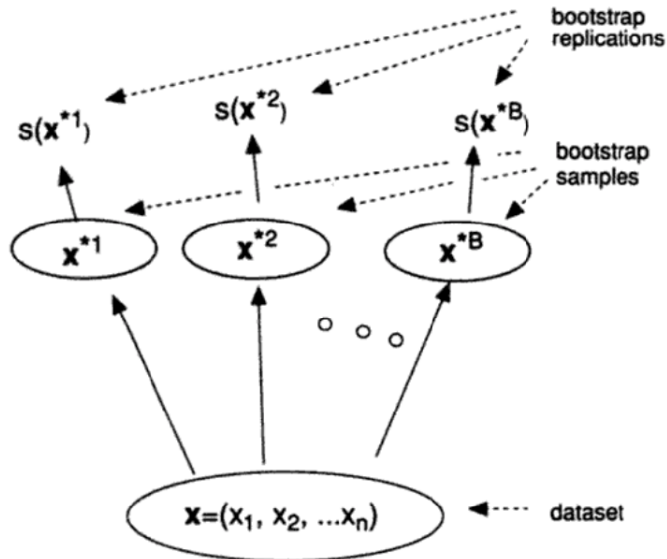
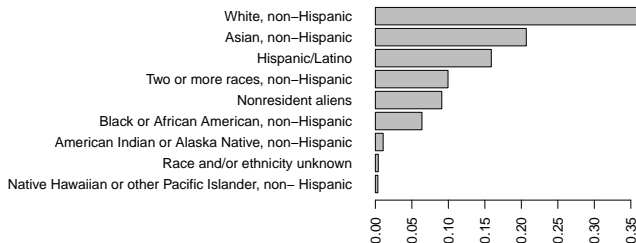


Illustration ¹

Données : origine ethnique des étudiants de Stanford



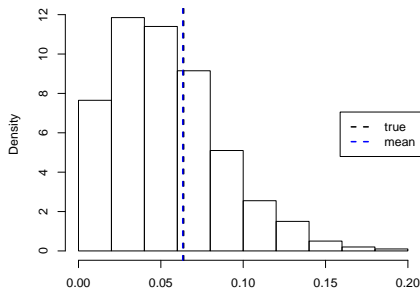
- ▶ On connaît toute la population
- ▶ Elle contient 6.4% d'étudiants afro-américains.

⇒ Question : estimer ce taux à partir d'un échantillon.

On considère :

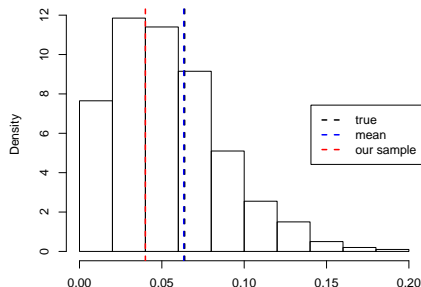
- ▶ des échantillons de taille 50.
- ▶ l'estimateur de la fréquence empirique.

⇒ **distribution d'échantillonnage** sur 1000 réalisations :

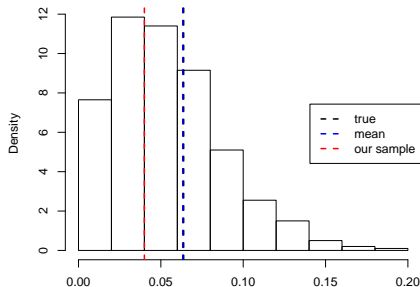


⇒ estimateur précis en moyenne, mais forte variance.

Sauf qu'en pratique, on n'aurait accès qu'à **une réalisation** :

[Introduction](#)[Formalisation](#)[Bootstrap pour l'inférence](#)[Caractérisation d'un estimateur](#)
[Intervalles de confiance](#)[Applications](#)[Bootstrap & régression](#)
[Bootstrap & prédiction](#)[Conclusion](#)[Références](#)

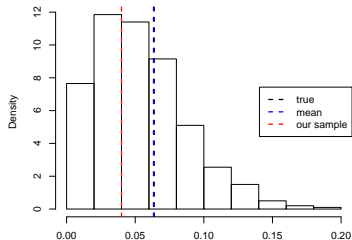
Sauf qu'en pratique, on n'aurait accès qu'à **une réalisation** :



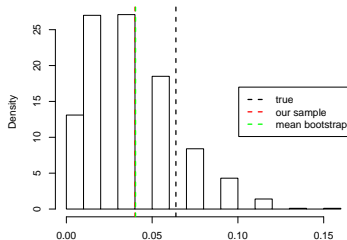
⇒ **bootstrap** : se baser uniquement sur notre échantillon :

- ▶ nombreux tirages avec remise
- ▶ distribution de la fréquence empirique "bootstrapée"

Tirages dans la **population** :



Tirages dans l'**échantillon** :



- ▶ l'échantillon \sim une population finie
- ▶ on simule des échantillons de cette population
- ▶ on veut en tirer des conclusions sur la vraie population

Remarques :

- ▶ intérêt limité sur cet exemple, car on connaît très bien les propriétés de l'estimateur de la fréquence empirique.
- ▶ intérêt général de la démarche :
 - ▶ ne pas "se forcer" à faire d'hypothèses quand les données ne s'y prêtent pas
 - ▶ permet de considérer des statistiques plus complexes
 - ▶ dont on ne connaît pas forcément la distribution d'échantillonnage

Bootstrap, vous avez dit bootstrap ?

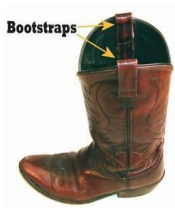
Terme introduit en statistique par **Bradley Efron** en 1979...

... vient de l'expression "**to pull oneself up by one's bootstrap**"...

- ▶ ~ s'en sortir par soi même, grâce à ses efforts

...souvent attribuée aux aventures du **Baron de Munchausen**.

- ▶ le Baron, tombé dans un marécage, s'en extrait en se tirant lui même par ses "bootstraps"



Formalisation du principe de ré-échantillonnage

Récapitulatif

On s'intéresse à une **variable aléatoire X** .

Outline

UE StatComp

Introduction

Formalisation

Bootstrap pour l'inférence

Caractérisation
d'un estimateur
Intervalle de
confiance

Applications

Bootstrap &
régression
Bootstrap &
prédiction

Conclusion

Références

Récapitulatif

On s'intéresse à une **variable aléatoire** X .

Au niveau de la **population** \mathcal{X} :

- ▶ X est régie par une **distribution** P et une **fonction de répartition** F :

$$P(X \leq x) = F(x), \quad \forall x \in \mathcal{X}.$$

- ▶ $\theta = t(P)$ est un **paramètre d'intérêt**

Récapitulatif

On s'intéresse à une **variable aléatoire** X .

Au niveau de la **population** \mathcal{X} :

- ▶ X est régie par une **distribution** P et une **fonction de répartition** F :

$$P(X \leq x) = F(x), \quad \forall x \in \mathcal{X}.$$

- ▶ $\theta = t(P)$ est un **paramètre d'intérêt**

On dispose d'un **échantillon** $\mathbf{X} = (X_1, \dots, X_n)$:

- ▶ les X_i sont iid selon P
- ▶ $\hat{\theta} = s(\mathbf{X})$ est un **estimateur** de θ

Récapitulatif

On s'intéresse à une **variable aléatoire** X .

Au niveau de la **population** \mathcal{X} :

- ▶ X est régie par une **distribution** P et une **fonction de répartition** F :

$$P(X \leq x) = F(x), \quad \forall x \in \mathcal{X}.$$

- ▶ $\theta = t(P)$ est un **paramètre d'intérêt**

On dispose d'un **échantillon** $\mathbf{X} = (X_1, \dots, X_n)$:

- ▶ les X_i sont iid selon P
- ▶ $\hat{\theta} = s(\mathbf{X})$ est un **estimateur** de θ

⇒ **inférence** : tirer des conclusions sur θ à partir de $\hat{\theta}$.

⇒ nécessite la **distribution d'échantillonnage** de $\hat{\theta}$.

Récapitulatif

Question clé : estimer la distribution d'échantillonnage de $\hat{\theta}$.

Outline

UE StatComp

Introduction

Formalisation

Bootstrap pour l'inférence

Caractérisation
d'un estimateur
Intervalle de
confiance

Applications

Bootstrap &
régression
Bootstrap &
prédiction

Conclusion

Références

Question clé : estimer la **distribution d'échantillonnage** de $\hat{\theta}$.

Stratégie #1 : l'estimer empiriquement à partir de **plusieurs échantillons** :

- ▶ on collecte plusieurs échantillons $\mathbf{X}_j = (X_{j1}, \dots, X_{jn})$
- ▶ on calcule $\hat{\theta}_j = s(\mathbf{X}_j)$
- ▶ on l'estime par la distribution des $\{\hat{\theta}_j\}$

Question clé : estimer la distribution d'échantillonnage de $\hat{\theta}$.

Stratégie #1 : l'estimer empiriquement à partir de plusieurs échantillons :

- ▶ on collecte plusieurs échantillons $\mathbf{X}_j = (X_{j1}, \dots, X_{jn})$
- ▶ on calcule $\hat{\theta}_j = s(\mathbf{X}_j)$
- ▶ on l'estime par la distribution des $\{\hat{\theta}_j\}$

Mais en pratique on ne dispose que d'un échantillon \mathbf{X} ...

Question clé : estimer la **distribution d'échantillonnage** de $\hat{\theta}$.

Stratégie #1 : l'estimer empiriquement à partir de **plusieurs échantillons** :

- ▶ on collecte plusieurs échantillons $\mathbf{X}_j = (X_{j1}, \dots, X_{jn})$
- ▶ on calcule $\hat{\theta}_j = s(\mathbf{X}_j)$
- ▶ on l'estime par la distribution des $\{\hat{\theta}_j\}$

Mais en pratique on ne dispose que d'un échantillon \mathbf{X} ...

⇒ **approche paramétrique** : faire des hypothèses sur la nature de la distribution P des données.

⇒ **approche non-paramétrique** : ré-échantillonnage dans \mathbf{X} pour estimer la distribution d'échantillonnage de $\hat{\theta}$.

Soit un échantillon $\mathbf{X} = (X_1, \dots, X_n)$.

On définit la **distribution empirique** \hat{P}_n comme :

$$\hat{P}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(X_i = x).$$

On définit de même la **fonction de répartition empirique** \hat{F}_n :

$$\hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(X_i \leq t).$$

Introduction

Formalisation

Bootstrap pour
l'inférenceCaractérisation
d'un estimateur
Intervalle de
confiance

Applications

Bootstrap &
régressionBootstrap &
prédiction

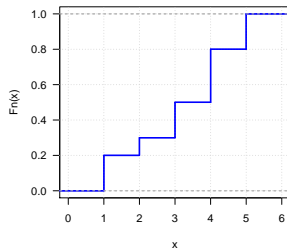
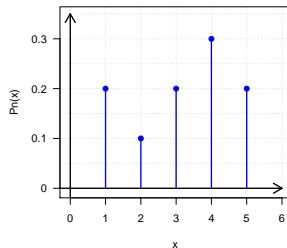
Conclusion

Références

Distribution empirique - illustration

Soit le vecteur $x = \{1, 1, 2, 3, 3, 4, 4, 4, 5, 5\}$.

\Rightarrow distribution empirique \hat{P}_n : \Rightarrow répartition empirique \hat{F}_n :

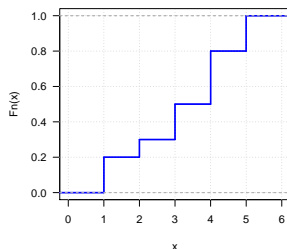
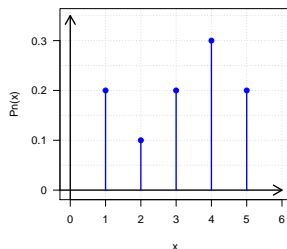


Distribution empirique & ré-échantillonnage

Outline

UE StatComp

Comment simuler un échantillon selon une **distribution empirique** ?



Introduction

Formalisation

Bootstrap pour l'inférence

Caractérisation d'un estimateur
Intervalles de confiance

Applications

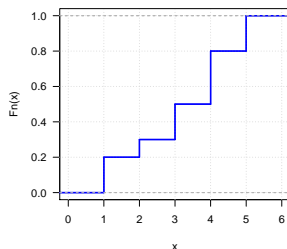
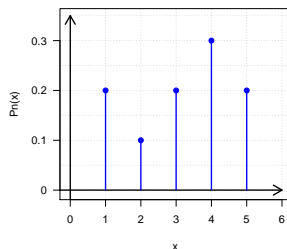
Bootstrap & régression
Bootstrap & prédiction

Conclusion

Références

Distribution empirique & ré-échantillonnage

Comment simuler un échantillon selon une **distribution empirique** ?



⇒ il suffit de **tirer avec remise** dans l'échantillon original.

⇒ pour s'en convaincre, on peut utiliser la méthode d'inversion.

Ré-échantillonnage & bootstrap

Outline

UE StatComp

Introduction

Formalisation

Bootstrap pour
l'inférence

Caractérisation
d'un estimateur
Intervalles de
confiance

Applications

Bootstrap &
régression

Bootstrap &
prédiction

Conclusion

Références

On s'intéresse à un estimateur $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$ d'un paramètre θ .

Ré-échantillonnage & bootstrap

Outline

UE StatComp

Introduction

Formalisation

Bootstrap pour
l'inférence

Caractérisation
d'un estimateur
Intervalle de
confiance

Applications

Bootstrap &
régression

Bootstrap &
prédiction

Conclusion

Références

On s'intéresse à un estimateur $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$ d'un paramètre θ .

On veut caractériser sa distribution d'échantillonnage, **en se basant uniquement sur l'échantillon (X_1, \dots, X_n) .**

On s'intéresse à un estimateur $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$ d'un paramètre θ .

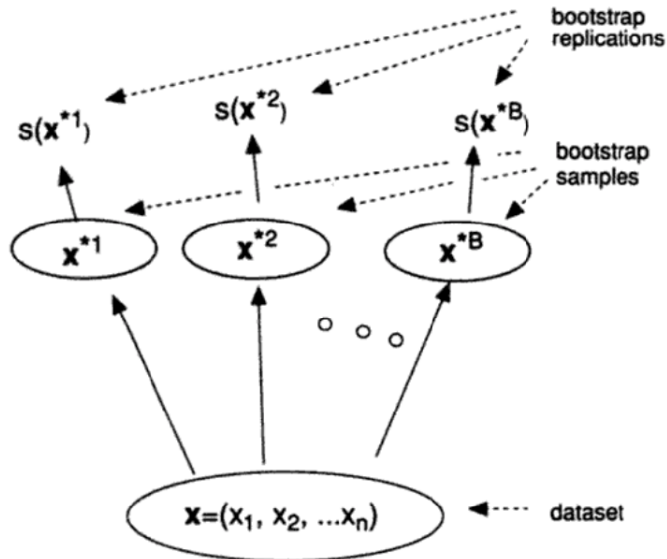
On veut caractériser sa distribution d'échantillonnage, **en se basant uniquement sur l'échantillon (X_1, \dots, X_n)** .

On applique la procédure suivante :

- ▶ Pour b allant de 1 à B ,
- ▶ On génère un échantillon (X_1^*, \dots, X_n^*) en **tirant avec remise** dans (X_1, \dots, X_n) .
- ▶ On calcule la valeur de notre estimateur $\hat{\theta}^{*(b)}$ à partir de (X_1^*, \dots, X_n^*) .

Et on travaille sur les B réalisations $(\hat{\theta}^{*(b)})_{b=1, \dots, B}$.

Ré-échantillonnage & bootstrap

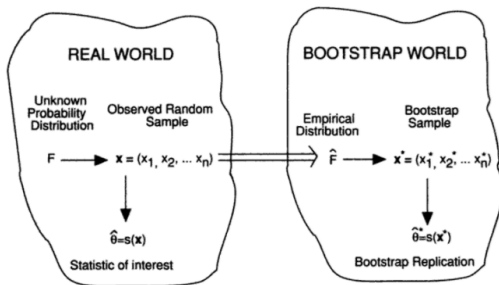


Monde réel et monde bootstrap

Formellement, le bootstrap considère des **réalisations tirées de la distribution empirique** définie par l'échantillon original.

On passe du **monde réel** au "monde bootstrap" :

- illustration tirée de Efron and Tibshirani (1993).



Monde réel et monde bootstrap

Monde réel vs monde bootstrap :

Outline

UE StatComp

Introduction

Formalisation

Bootstrap pour
l'inférence

Caractérisation
d'un estimateur
Intervalle de
confiance

Applications

Bootstrap &
régression

Bootstrap &
prédiction

Conclusion

Références

Monde réel et monde bootstrap

Outline

UE StatComp

Introduction

Formalisation

Bootstrap pour
l'inférence

Caractérisation
d'un estimateur
Intervalle de
confiance

Applications

Bootstrap &
régression

Bootstrap &
prédiction

Conclusion

Références

Monde réel vs monde bootstrap :

- ▶ distribution des données : P vs \hat{P}_n
 - ▶ vraie distribution P inconnue
 - ▶ distribution \hat{P}_n parfaitement connue

Monde réel vs monde bootstrap :

- ▶ **distribution des données** : P vs \hat{P}_n
 - ▶ vraie distribution P inconnue
 - ▶ distribution \hat{P}_n parfaitement connue
- ▶ **échantillon** : (X_1, \dots, X_n) vs (X_1^*, \dots, X_n^*)
 - ▶ un unique échantillon $(X_1, \dots, X_n) \sim P$
 - ▶ autant d'échantillons $(X_1^*, \dots, X_n^*) \sim \hat{P}_n$ qu'on veut

Monde réel vs monde bootstrap :

- ▶ **distribution des données** : P vs \hat{P}_n
 - ▶ vraie distribution P inconnue
 - ▶ distribution \hat{P}_n parfaitement connue
- ▶ **échantillon** : (X_1, \dots, X_n) vs (X_1^*, \dots, X_n^*)
 - ▶ un unique échantillon $(X_1, \dots, X_n) \sim P$
 - ▶ autant d'échantillons $(X_1^*, \dots, X_n^*) \sim \hat{P}_n$ qu'on veut
- ▶ **paramètre** : θ vs $\hat{\theta}$ vs $\hat{\theta}^*$
 - ▶ un vrai paramètre θ inconnu
 - ▶ une estimation $\hat{\theta} = s(X_1, \dots, X_n)$ connue
 - ▶ NB : une estimation de θ
 - ▶ autant d'estimations $\hat{\theta}^* = s(X_1^*, \dots, X_n^*)$ qu'on veut
 - ▶ NB : des estimations de $\hat{\theta}$

Introduction

Formalisation

Bootstrap pour
l'inférenceCaractérisation
d'un estimateur
Intervalles de
confiance

Applications

Bootstrap &
régressionBootstrap &
prédiction

Conclusion

Références

Monde réel et monde bootstrap

Outline

UE StatComp

Introduction

Formalisation

Bootstrap pour
l'inférence

Caractérisation
d'un estimateur
Intervalle de
confiance

Applications

Bootstrap &
régression
Bootstrap &
prédiction

Conclusion

Références

Dans le **monde bootstrap** :

- ▶ la distribution des données est \hat{P}_n
 - ▶ elle est **parfaitement connue**
 - ▶ on peut en **tirer/simuler des échantillons**
- ▶ le "vrai" paramètre de la population est $\hat{\theta}$
 - ▶ il est **parfaitement connu**
 - ▶ on peut le **comparer aux estimations** $\hat{\theta}^*$

Dans le **monde bootstrap** :

- ▶ la distribution des données est \hat{P}_n
 - ▶ elle est **parfaitement connue**
 - ▶ on peut en **tirer/simuler des échantillons**
- ▶ le "vrai" paramètre de la population est $\hat{\theta}$
 - ▶ il est **parfaitement connu**
 - ▶ on peut le **comparer aux estimations** $\hat{\theta}^*$

⇒ **principe du bootstrap** :

1. se placer dans le monde bootstrap
2. calculer la distribution d'échantillonnage des $\hat{\theta}^*$
3. se comparer à $\hat{\theta}$ (le "vrai" paramètre de \hat{P}_n)
4. en déduire des caractéristiques de $\hat{\theta}$ (par rapport à P).
 - ▶ e.g., biais, erreur-type et intervalle de confiance

Introduction

Formalisation

Bootstrap pour l'inférence

Caractérisation d'un estimateur
Intervalles de confiance

Applications

Bootstrap & régression
Bootstrap & prédiction

Conclusion

Références

Bootstrap et inférence statistique : caractérisation d'un estimateur

On peut par exemple appliquer le bootstrap pour :

- ▶ Estimer le **biais** d'un estimateur.
- ▶ Estimer son **erreur quadratique moyenne**.
- ▶ Estimer son **erreur type**.
- ▶ Donner un intervalle de confiance sur une estimation.

Estimer le biais d'un estimateur

Rappel : biais d'un estimateur : $\text{Biais}(\hat{\theta}) = E[\hat{\theta}] - \theta$.

Outline

UE StatComp

Introduction

Formalisation

Bootstrap pour l'inférence

Caractérisation d'un estimateur

Intervalle de confiance

Applications

Bootstrap & régression

Bootstrap & prédiction

Conclusion

Références

Estimer le biais d'un estimateur

Rappel : biais d'un estimateur : $\text{Biais}(\hat{\theta}) = E[\hat{\theta}] - \theta$.

Procédure bootstrap :

1. On dispose d'un échantillon (X_1, \dots, X_n) .
2. On calcule $\hat{\theta}$ sur l'échantillon.
3. On applique la procédure bootstrap pour obtenir des réalisations $\hat{\theta}^{*(b)}$:
 - ▶ Pour b allant de 1 à B ,
 - ▶ on génère un échantillon (X_1^*, \dots, X_n^*) ,
 - ▶ on calcule $\hat{\theta}^{*(b)}$ à partir de (X_1^*, \dots, X_n^*) .

Estimer le biais d'un estimateur

Rappel : **biais** d'un estimateur : $\text{Biais}(\hat{\theta}) = E[\hat{\theta}] - \theta$.

Procédure bootstrap :

1. On dispose d'un échantillon (X_1, \dots, X_n) .
2. On calcule $\hat{\theta}$ sur l'échantillon.
3. On applique la procédure bootstrap pour obtenir des réalisations $\hat{\theta}^{*(b)}$:
 - ▶ Pour b allant de 1 à B ,
 - ▶ on génère un échantillon (X_1^*, \dots, X_n^*) ,
 - ▶ on calcule $\hat{\theta}^{*(b)}$ à partir de (X_1^*, \dots, X_n^*) .

On estime le biais de $\hat{\theta}$ par :

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

Estimer l'erreur quadratique moyenne (MSE) d'un estimateur

Rappel : MSE d'un estimateur : $\text{MSE}(\hat{\theta}) = E[(\hat{\theta} - \theta)^2]$.

Estimer l'erreur quadratique moyenne (MSE) d'un estimateur

Rappel : MSE d'un estimateur : $\text{MSE}(\hat{\theta}) = E[(\hat{\theta} - \theta)^2]$.

Procédure bootstrap :

1. On dispose d'un échantillon (X_1, \dots, X_n) .
2. On calcule $\hat{\theta}$ sur l'échantillon.
3. On applique la procédure bootstrap pour obtenir des réalisations $\hat{\theta}^{*(b)}$:
 - ▶ Pour b allant de 1 à B ,
 - ▶ on génère un échantillon (X_1^*, \dots, X_n^*) ,
 - ▶ on calcule $\hat{\theta}^{*(b)}$ à partir de (X_1^*, \dots, X_n^*) .

Estimer l'erreur quadratique moyenne (MSE) d'un estimateur

Rappel : MSE d'un estimateur : $\text{MSE}(\hat{\theta}) = E[(\hat{\theta} - \theta)^2]$.

Procédure bootstrap :

1. On dispose d'un échantillon (X_1, \dots, X_n) .
2. On calcule $\hat{\theta}$ sur l'échantillon.
3. On applique la procédure bootstrap pour obtenir des réalisations $\hat{\theta}^{*(b)}$:
 - ▶ Pour b allant de 1 à B ,
 - ▶ on génère un échantillon (X_1^*, \dots, X_n^*) ,
 - ▶ on calcule $\hat{\theta}^{*(b)}$ à partir de (X_1^*, \dots, X_n^*) .

On estime l'erreur quadratique moyenne de $\hat{\theta}$ par :

$$\widehat{\text{MSE}}(\hat{\theta}) = \frac{1}{B} \sum_{b=1}^B (\hat{\theta}^{*(b)} - \hat{\theta})^2.$$

Estimer l'erreur type d'un estimateur

Rappel : erreur type d'un estimateur : l'écart type de sa distribution d'échantillonnage.

Outline

UE StatComp

Introduction

Formalisation

Bootstrap pour
l'inférence

**Caractérisation
d'un estimateur**
Intervalles de
confiance

Applications

Bootstrap &
régression

Bootstrap &
prédiction

Conclusion

Références

Estimer l'erreur type d'un estimateur

Rappel : **erreur type** d'un estimateur : l'écart type de sa distribution d'échantillonnage.

Procédure bootstrap :

1. On dispose d'un échantillon (X_1, \dots, X_n) .
2. On applique la procédure bootstrap pour obtenir des réalisations $\hat{\theta}^{*(b)}$:
 - ▶ Pour b allant de 1 à B ,
 - ▶ on génère un échantillon (X_1^*, \dots, X_n^*) ,
 - ▶ on calcule $\hat{\theta}^{*(b)}$ à partir de (X_1^*, \dots, X_n^*) .

Estimer l'erreur type d'un estimateur

Rappel : **erreur type** d'un estimateur : l'écart type de sa distribution d'échantillonnage.

Procédure bootstrap :

1. On dispose d'un échantillon (X_1, \dots, X_n) .
2. On applique la procédure bootstrap pour obtenir des réalisations $\hat{\theta}^{*(b)}$:
 - ▶ Pour b allant de 1 à B ,
 - ▶ on génère un échantillon (X_1^*, \dots, X_n^*) ,
 - ▶ on calcule $\hat{\theta}^{*(b)}$ à partir de (X_1^*, \dots, X_n^*) .

On estime l'erreur-type de $\hat{\theta}$ par l'écart-type des $(\hat{\theta}^{*(b)})$:

$$\widehat{\text{se}}(\hat{\theta}) = \sqrt{\frac{1}{B-1} \sum_{b=1}^B \left(\hat{\theta}^{*(b)} - \bar{\hat{\theta}}^* \right)^2} \quad \text{où} \quad \bar{\hat{\theta}}^* = \frac{1}{B} \sum_{b=1}^B \hat{\theta}^{*(b)}.$$

Remarque importante

In fine, on estime la "vraie" variance de $\hat{\theta}$ (i.e., selon P) par la variance empirique des $(\hat{\theta}^{*(b)})$ (i.e., selon \hat{P}_n).

Remarque importante

In fine, on estime la "vraie" variance de $\hat{\theta}$ (i.e., selon P) par la variance empirique des $(\hat{\theta}^{*(b)})$ (i.e., selon \hat{P}_n).

Il faut bien comprendre qu'il y a 2 niveaux d'approximation :

1. approximer la vraie distribution P par l'empirique \hat{P}_n .
2. approximer la variance de $\hat{\theta}$ selon \hat{P}_n par la variance empirique des $(\hat{\theta}^{*(b)})$.

Remarque importante

In fine, on estime la "vraie" variance de $\hat{\theta}$ (i.e., selon P) par la variance empirique des $(\hat{\theta}^{*(b)})$ (i.e., selon \hat{P}_n).

Il faut bien comprendre qu'il y a 2 niveaux d'approximation :

1. approximer la vraie distribution P par l'empirique \hat{P}_n .
2. approximer la variance de $\hat{\theta}$ selon \hat{P}_n par la variance empirique des $(\hat{\theta}^{*(b)})$.

⇒ point #2 : ok, il suffit de prendre B grand.

⇒ point #1 : plus délicat...mais valide quand n est grand.

Remarque importante

In fine, on estime la "vraie" variance de $\hat{\theta}$ (i.e., selon P) par la variance empirique des $(\hat{\theta}^{*(b)})$ (i.e., selon \hat{P}_n).

Il faut bien comprendre qu'il y a 2 niveaux d'approximation :

1. approximer la vraie distribution P par l'empirique \hat{P}_n .
2. approximer la variance de $\hat{\theta}$ selon \hat{P}_n par la variance empirique des $(\hat{\theta}^{*(b)})$.

⇒ point #2 : ok, il suffit de prendre B grand.

⇒ point #1 : plus délicat...mais valide quand n est grand.

⚠ bootstrap \neq méthode pour petits échantillons !

- ▶ intérêt #1 = relâcher hypothèses paramétriques
- ▶ intérêt #2 = générique, valable pour toute statistique

Calcul du biais et de l'erreur type de la moyenne :

```
> x = c(1,1,2,3,3,4,4,4,5,5)
> n = length(x)
> B = 2000
> theta.hat = mean(x)
> theta.b = numeric(B)
> for(b in 1:B){
  ind = sample(1:n, size = n, replace = TRUE)
  theta.b[b] = mean(x[ind])
}
> bias = mean(theta.b) - theta.hat
> se = sd(theta.b)
```

Bootstrap et inférence statistique : intervalles de confiance

Bootstrap et intervalles de confiance

Il existe de nombreuses manières de définir des intervalles de confiance par bootstrap.

Outline

UE StatComp

Introduction

Formalisation

Bootstrap pour l'inférence

Caractérisation d'un estimateur

Intervalles de confiance

Applications

Bootstrap & régression

Bootstrap & prédiction

Conclusion

Références

Il existe de **nombreuses manières** de définir des intervalles de confiance par bootstrap.

Nous allons considérer deux définitions se basant uniquement sur les quantiles de la distribution des $\hat{\theta}^{*(b)}$:

- ▶ l'intervalle de confiance **des percentiles**.
- ▶ l'intervalle de confiance **basique** (ou "du pivot").

Ces définitions font le moins d'hypothèses possible.

Intervalle de confiance - méthode des percentiles

Cette définition est probablement la plus simple et intuitive.

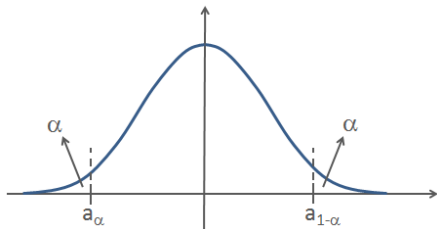
Elle consiste à calculer empiriquement l'**intervalle couvrant**
 $100(1 - \alpha)\%$ des estimations bootstrap obtenues.

Formellement, en notant q_{α}^* le **quantile d'ordre α de la**
distribution des estimations bootstrap $\hat{\theta}^{*(b)}$, il est défini
comme :

$$[q_{\alpha/2}^* ; q_{1-\alpha/2}^*].$$

Intervalle de confiance - méthode basique (1/3)

On va s'intéresser à la distribution de la statistique $(\hat{\theta} - \theta)$:



Si a_α le quantile d'ordre α de la statistique $(\hat{\theta} - \theta)$ alors :

$$P(\hat{\theta} - \theta \leq a_\alpha) = \alpha \text{ et } P(\hat{\theta} - \theta \geq a_{1-\alpha}) = \alpha.$$

Intervalle de confiance - méthode basique (2/3)

On a donc (par définition) :

$$P(\hat{\theta} - \theta \leq a_{\alpha}) = \alpha \text{ et } P(\hat{\theta} - \theta \geq a_{1-\alpha}) = \alpha,$$

où a_{α} le quantile d'ordre α de la statistique $(\hat{\theta} - \theta)$.

Intervalle de confiance - méthode basique (2/3)

On a donc (par définition) :

$$P(\hat{\theta} - \theta \leq a_{\alpha}) = \alpha \text{ et } P(\hat{\theta} - \theta \geq a_{1-\alpha}) = \alpha,$$

où a_{α} le quantile d'ordre α de la statistique $(\hat{\theta} - \theta)$.

\Rightarrow on peut écrire :

$$P(\theta \geq \hat{\theta} - a_{\alpha}) = \alpha \text{ et } P(\theta \leq \hat{\theta} - a_{1-\alpha}) = \alpha,$$

Intervalle de confiance - méthode basique (2/3)

On a donc (par définition) :

$$P(\hat{\theta} - \theta \leq a_{\alpha}) = \alpha \text{ et } P(\hat{\theta} - \theta \geq a_{1-\alpha}) = \alpha,$$

où a_{α} le quantile d'ordre α de la statistique $(\hat{\theta} - \theta)$.

\Rightarrow on peut écrire :

$$P(\theta \geq \hat{\theta} - a_{\alpha}) = \alpha \text{ et } P(\theta \leq \hat{\theta} - a_{1-\alpha}) = \alpha,$$

et en déduire l'intervalle de confiance à $100(1 - \alpha)\%$:

$$[\hat{\theta} - a_{1-\alpha/2} ; \hat{\theta} - a_{\alpha/2}].$$

Intervalle de confiance - méthode basique (3/3)

On considère donc l'intervalle de confiance défini comme

$$[\hat{\theta} - a_{1-\alpha/2} ; \hat{\theta} - a_{\alpha/2}],$$

où a_{α} le quantile d'ordre α de la statistique $(\hat{\theta} - \theta)$.

Intervalle de confiance - méthode basique (3/3)

On considère donc l'intervalle de confiance défini comme

$$[\hat{\theta} - a_{1-\alpha/2} ; \hat{\theta} - a_{\alpha/2}],$$

où a_{α} le quantile d'ordre α de la statistique $(\hat{\theta} - \theta)$.

Problème : on ne connaît pas la distribution de $(\hat{\theta} - \theta)$.

Intervalle de confiance - méthode basique (3/3)

On considère donc l'intervalle de confiance défini comme

$$[\hat{\theta} - a_{1-\alpha/2} ; \hat{\theta} - a_{\alpha/2}],$$

où a_{α} le quantile d'ordre α de la statistique $(\hat{\theta} - \theta)$.

Problème : on ne connaît pas la distribution de $(\hat{\theta} - \theta)$.

⇒ On passe dans le **monde bootstrap** : on approxime la distribution de $(\hat{\theta} - \theta)$ par celle de $(\hat{\theta}^{*(b)} - \hat{\theta})$.

Intervalle de confiance - méthode basique (3/3)

On considère donc l'intervalle de confiance défini comme

$$[\hat{\theta} - a_{1-\alpha/2} ; \hat{\theta} - a_{\alpha/2}],$$

où a_α le quantile d'ordre α de la statistique $(\hat{\theta} - \theta)$.

Problème : on ne connaît pas la distribution de $(\hat{\theta} - \theta)$.

⇒ On passe dans le **monde bootstrap** : on approxime la distribution de $(\hat{\theta} - \theta)$ par celle de $(\hat{\theta}^{*(b)} - \hat{\theta})$.

► on peut donc écrire :

$$a_\alpha = q_\alpha^* - \hat{\theta}, \text{ où } q_\alpha^* \text{ est le quantile d'ordre } \alpha \text{ des } \hat{\theta}^{*(b)}.$$

Intervalle de confiance - méthode basique (3/3)

On considère donc l'intervalle de confiance défini comme

$$[\hat{\theta} - a_{1-\alpha/2} ; \hat{\theta} - a_{\alpha/2}],$$

où a_α le quantile d'ordre α de la statistique $(\hat{\theta} - \theta)$.

Problème : on ne connaît pas la distribution de $(\hat{\theta} - \theta)$.

⇒ On passe dans le **monde bootstrap** : on approxime la distribution de $(\hat{\theta} - \theta)$ par celle de $(\hat{\theta}^{*(b)} - \hat{\theta})$.

► on peut donc écrire :

$$a_\alpha = q_\alpha^* - \hat{\theta}, \text{ où } q_\alpha^* \text{ est le quantile d'ordre } \alpha \text{ des } \hat{\theta}^{*(b)}.$$

► on en déduit la définition suivante :

$$[2\hat{\theta} - q_{1-\alpha/2}^* ; 2\hat{\theta} - q_{\alpha/2}^*].$$

Bootstrap & intervalles de confiance - remarques

Deux définitions considérées :

1. IC des percentiles : $[q_{\alpha/2}^* ; q_{1-\alpha/2}^*]$.
2. IC basique : $[2\hat{\theta} - q_{1-\alpha/2}^* ; 2\hat{\theta} - q_{\alpha/2}^*]$.

\Rightarrow se basent uniquement sur q_{α}^* : le quantile d'ordre α des estimations bootstrap $\hat{\theta}^{*(b)}$.

\Rightarrow méthode "basique" / "pivot" : analogie monde bootstrap

► $\theta \rightarrow \hat{\theta} ; \hat{\theta} \rightarrow \hat{\theta}^*$

Bootstrap & intervalles de confiance - remarques

Deux définitions considérées :

1. IC des percentiles : $[q_{\alpha/2}^* ; q_{1-\alpha/2}^*]$.
2. IC basique : $[2\hat{\theta} - q_{1-\alpha/2}^* ; 2\hat{\theta} - q_{\alpha/2}^*]$.

⇒ se basent uniquement sur q_{α}^* : le quantile d'ordre α des estimations bootstrap $\hat{\theta}^{*(b)}$.

⇒ méthode "basique" / "pivot" : analogie monde bootstrap

$$\blacktriangleright \theta \rightarrow \hat{\theta} ; \hat{\theta} \rightarrow \hat{\theta}^*$$

De nombreuses autres définitions existent.

- ▶ normal, "studentisé", accéléré, corrigé du biais...

Applications : bootstrap et régression

Objectif : prédire / modéliser une variable $Y \in \mathbb{R}$ à partir de p variables explicatives $X^j \in \mathbb{R}$.

Modèle :

$$Y = \beta_0 + \sum_{j=1}^p \beta_j X^j + \epsilon \quad \text{avec } \epsilon \text{ un terme d'erreur.}$$

Introduction

Formalisation

Bootstrap pour
l'inférence

Caractérisation
d'un estimateur
Intervalle de
confiance

Applications

**Bootstrap &
régression**

Bootstrap &
prédiction

Conclusion

Références

Objectif : prédire / modéliser une variable $Y \in \mathbb{R}$ à partir de p variables explicatives $X^j \in \mathbb{R}$.

Modèle :

$$Y = \beta_0 + \sum_{j=1}^p \beta_j X^j + \epsilon \quad \text{avec } \epsilon \text{ un terme d'erreur.}$$

\Rightarrow on estime les coefficients β_j par **moindre carrés**

- ▶ à partir d'un échantillon $(X_i, Y_i) \in \mathbb{R}^p \times \mathbb{R}, i = 1, \dots, n$

Introduction

Formalisation

Bootstrap pour
l'inférenceCaractérisation
d'un estimateur
Intervalle de
confiance

Applications

**Bootstrap &
régression**Bootstrap &
prédiction

Conclusion

Références

Objectif : prédire / modéliser une variable $Y \in \mathbb{R}$ à partir de p variables explicatives $X^j \in \mathbb{R}$.

Modèle :

$$Y = \beta_0 + \sum_{j=1}^p \beta_j X^j + \epsilon \quad \text{avec } \epsilon \text{ un terme d'erreur.}$$

\Rightarrow on estime les coefficients β_j par **moindre carrés**

▶ à partir d'un échantillon $(X_i, Y_i) \in \mathbb{R}^p \times \mathbb{R}$, $i = 1, \dots, n$

\Rightarrow **sous l'hypothèse que les résidus ϵ_i sont iid selon $\mathcal{N}(0, \sigma^2)$**
on connaît la distribution d'échantillonnage des $\hat{\beta}_j$.

▶ on peut donc en tirer des intervalles de confiance

Bootstrap et régression = alternative non-paramétrique.

- ▶ relâcher les hypothèses du modèle.

Objectif : estimer la distribution d'échantillonnage des $\hat{\beta}_j$.

- ▶ in fine : calcul d'intervalles de confiance comme avant.

Deux stratégies principales :

1. bootstrap **par paires**
2. bootstrap **des résidus**

Principe : tirer avec remise dans $\{(X_i, Y_i)\}_{i=1, \dots, n}$.

Procédure bootstrap :

- ▶ On dispose d'un échantillon (Z_1, \dots, Z_n) , $Z_i = (X_i, Y_i)$
- ▶ On estime $\hat{\beta}_j$ sur l'échantillon original
- ▶ On applique le **bootstrap par paires** :
 - ▶ pour b de 1 à B
 - ▶ on génère un échantillon (Z_1^*, \dots, Z_n^*)
 - ▶ on estime les coefficients $\hat{\beta}_j^*$ à partir des (Z_1^*, \dots, Z_n^*)

Principe : tirer avec remise dans $\{(X_i, Y_i)\}_{i=1, \dots, n}$.

Procédure bootstrap :

- ▶ On dispose d'un échantillon (Z_1, \dots, Z_n) , $Z_i = (X_i, Y_i)$
- ▶ On estime $\hat{\beta}_j$ sur l'échantillon original
- ▶ On applique le **bootstrap par paires** :
 - ▶ pour b de 1 à B
 - ▶ on génère un échantillon (Z_1^*, \dots, Z_n^*)
 - ▶ on estime les coefficients $\hat{\beta}_j^*$ à partir des (Z_1^*, \dots, Z_n^*)

\Rightarrow l'approche "standard".

\Rightarrow hypothèse : (X_i, Y_i) iid selon une loi (jointe) P .

Bootstrap des résidus

Principe : travailler à partir des **résidus** du modèle initial.

Procédure bootstrap :

- ▶ On dispose d'un échantillon (Z_1, \dots, Z_n) , $Z_i = (X_i, Y_i)$
- ▶ On estime $\hat{\beta}_j$ sur l'échantillon original
- ▶ On considère les résidus $(\epsilon_1, \dots, \epsilon_n)$
- ▶ On applique le **bootstrap par résidus** :
 - ▶ pour b de 1 à B
 - ▶ on génère un échantillon $(\epsilon_1^*, \dots, \epsilon_n^*)$
 - ▶ on calcule $Y_i^* = \hat{\beta}_0 + \sum_{j=1}^p \hat{\beta}_j X_{ij} + \epsilon_i^*$
 - ▶ on estime les coefficients $\hat{\beta}_j^*$ à partir des (X_i, Y_i^*)

Bootstrap des résidus

Principe : travailler à partir des **résidus** du modèle initial.

Procédure bootstrap :

- ▶ On dispose d'un échantillon (Z_1, \dots, Z_n) , $Z_i = (X_i, Y_i)$
- ▶ On estime $\hat{\beta}_j$ sur l'échantillon original
- ▶ On considère les résidus $(\epsilon_1, \dots, \epsilon_n)$
- ▶ On applique le **bootstrap par résidus** :
 - ▶ pour b de 1 à B
 - ▶ on génère un échantillon $(\epsilon_1^*, \dots, \epsilon_n^*)$
 - ▶ on calcule $Y_i^* = \hat{\beta}_0 + \sum_{j=1}^p \hat{\beta}_j X_{ij} + \epsilon_i^*$
 - ▶ on estime les coefficients $\hat{\beta}_j^*$ à partir des (X_i, Y_i^*)

⇒ on tire avec remise les **résidus**.

- ▶ tous les X_i sont utilisés à chaque fois
- ▶ hypothèse : $Y_i|X_i$ iid

⇒ semble être la stratégie la plus classique

Applications : bootstrap et modèles de prédiction

Stratégie bootstrap & "bagging"

Principe bootstrap pour la construction de prédicteurs :

Stratégie bootstrap & "bagging"

Principe bootstrap pour la construction de prédicteurs :

1. **Apprentissage** : construire B prédicteurs à partir d'échantillons obtenus en tirant avec remise dans l'échantillon original.

Introduction

Formalisation

Bootstrap pour l'inférence

Caractérisation d'un estimateur
Intervalle de confiance

Applications

Bootstrap & régression
Bootstrap & prédiction

Conclusion

Références

Stratégie bootstrap & "bagging"

Principe bootstrap pour la construction de prédicteurs :

1. **Apprentissage** : construire B prédicteurs à partir d'échantillons obtenus en tirant avec remise dans l'échantillon original.
2. **Prédiction** : agréger les prédictions des B modèles
 - ▶ **régression** : prédire la moyenne des valeurs obtenues.
 - ▶ **classification** : prédire la classe prédite le plus souvent

Stratégie bootstrap & "bagging"

Principe bootstrap pour la construction de prédicteurs :

1. **Apprentissage** : construire B prédicteurs à partir d'échantillons obtenus en tirant avec remise dans l'échantillon original.
2. **Prédiction** : agréger les prédictions des B modèles
 - ▶ **régression** : prédire la moyenne des valeurs obtenues.
 - ▶ **classification** : prédire la classe prédite le plus souvent

On parle de stratégie **bagging**, pour **bootstrap-aggregating**.

Stratégie générique, souvent basée sur des arbres de décision

En pratique elle permet de **limiter le sur-apprentissage**.

Stratégie bootstrap & "bagging"

Principe bootstrap pour la construction de prédicteurs :

1. **Apprentissage** : construire B prédicteurs à partir d'échantillons obtenus en tirant avec remise dans l'échantillon original.
2. **Prédiction** : agréger les prédictions des B modèles
 - ▶ **régression** : prédire la moyenne des valeurs obtenues.
 - ▶ **classification** : prédire la classe prédite le plus souvent

On parle de stratégie **bagging**, pour **bootstrap-aggregating**.

Stratégie générique, souvent basée sur des arbres de décision

En pratique elle permet de **limiter le sur-apprentissage**.

⇒ à suivre dans cours **"Apprentissage Statistique II"**.

Conclusion

- Bootstrap : un principe très simple à mettre en oeuvre.

- ▶ Bootstrap : un principe très simple à mettre en oeuvre.
- ▶ Il permet de répondre à des questions d'inférence statistique sans aucune information sur la loi de la variable aléatoire sous-jacente.

- ▶ Bootstrap : un principe très simple à mettre en oeuvre.
- ▶ Il permet de répondre à des questions d'inférence statistique sans aucune information sur la loi de la variable aléatoire sous-jacente.
- ▶ Dans ce cas là, travailler par ré-échantillonnage de l'échantillon disponible est parfois la meilleure stratégie, surtout si la loi sous-jacente n'est pas une loi usuelle.

► Intérêts principaux :

1. relâcher les hypothèses sur la loi de la variable aléatoire étudiée qu'on doit faire avec les approches paramétriques
2. principe générique applicable à n'importe quelle statistique (dont on ne connaît pas la distribution)

► Intérêts principaux :

1. relâcher les hypothèses sur la loi de la variable aléatoire étudiée qu'on doit faire avec les approches paramétriques
2. principe générique applicable à n'importe quelle statistique (dont on ne connaît pas la distribution)

► Bien garder en tête qu'il y a 2 niveaux d'approximation

1. remplacer la "vraie" distribution par l'empirique
2. remplacer la "vraie" distribution d'échantillonnage (selon la loi empirique) par celle obtenue par tirages

► Intérêts principaux :

1. relâcher les hypothèses sur la loi de la variable aléatoire étudiée qu'on doit faire avec les approches paramétriques
2. principe générique applicable à n'importe quelle statistique (dont on ne connaît pas la distribution)

► Bien garder en tête qu'il y a 2 niveaux d'approximation

1. remplacer la "vraie" distribution par l'empirique
2. remplacer la "vraie" distribution d'échantillonnage (selon la loi empirique) par celle obtenue par tirages

► Par conséquent : méthode valide quand n est grand

- résultats théoriques pour démontrer la validité des procédures décrites (caractérisation et IC)
- le bootstrap n'est **pas** dédié aux petits échantillons

► Intérêts principaux :

1. relâcher les hypothèses sur la loi de la variable aléatoire étudiée qu'on doit faire avec les approches paramétriques
2. principe générique applicable à n'importe quelle statistique (dont on ne connaît pas la distribution)

► Bien garder en tête qu'il y a 2 niveaux d'approximation

1. remplacer la "vraie" distribution par l'empirique
2. remplacer la "vraie" distribution d'échantillonnage (selon la loi empirique) par celle obtenue par tirages

► Par conséquent : méthode valide quand n est grand

- résultats théoriques pour démontrer la validité des procédures décrites (caractérisation et IC)
- le bootstrap n'est **pas** dédié aux petits échantillons

► Et si n est petit ? Pas forcément pire qu'une approche paramétrique...

- Le principe du bootstrap a été décliné avec succès pour **construire des modèles de prédiction**, dans une stratégie dite de **bagging**.

- ▶ Le principe du bootstrap a été décliné avec succès pour **construire des modèles de prédiction**, dans une stratégie dite de **bagging**.
- ▶ Un exemple important est celui des **forêts aléatoires**, qui sont des **classifieurs très performants** et relativement simples à mettre en oeuvre.

- ▶ Le principe du bootstrap a été décliné avec succès pour **construire des modèles de prédiction**, dans une stratégie dite de **bagging**.
- ▶ Un exemple important est celui des **forêts aléatoires**, qui sont des **classifieurs très performants** et relativement simples à mettre en oeuvre.
- ▶ On peut également utiliser ce principe pour **évaluer les performances d'un modèle de prédiction**, comme une alternative à la **validation croisée**.
 - ▶ mais c'est moins classique.

⇒ à suivre dans le cours "**Apprentissage Statistique II**".

Le pilier du bootstrap :

```
> ind.bs = sample(n, replace = TRUE)
```

Les packages **bootstrap** et **boot** implémentent les méthodes vues en cours.

Le package **boot** est le plus recommandé.

- ▶ extrait tiré de la documentation du package **bootstrap** : *New projects should preferentially use the recommended package "boot"*

Package `boot` : deux fonctions principales

1. fonction `boot()`

- ▶ en entrée : données, statistique considérée (une fonction) et B .
- ▶ en sortie : un objet contenant (en 1er lieu) les estimations bootstrap.
- ▶ la fonction `print()` affiche le biais et l'erreur type.

2. fonction `boot.ci()`

- ▶ en entrée : l'objet retourné par la fonction `boot()`.
- ▶ en sortie : les intervalles de confiance par différentes stratégies (e.g., `basic`, `perc`, `norm`).

⇒ à voir en TP.

B. Efron and R. Tibshirani. *An introduction to the bootstrap*.
Chapman & Hall, 1993.