

**ÉCOLE DU PÉTROLE ET DES MOTEURS**

**INSTITUT FRANÇAIS DU PÉTROLE**

228-232, avenue Napoléon Bonaparte

92852 RUEIL-MALMAISON CEDEX

téléphone : 01 47 52 62 80 - télécopieur : 01 47 52 70 36

**La mise en oeuvre des techniques de Bootstrap  
pour la prévision économétrique :  
application à l'industrie automobile**

*Sandrine JUAN\**

*Frédéric LANTZ*

novembre 2000

**Cahiers de l'économie - n° 39**

**Série Recherche**

\*Renault, 1, avenue du Golf, 78288 GUYANCOURT CEDEX – sandrine.juan@renault.com

La collection "Les cahiers de l'économie" a pour objectif de présenter des travaux réalisés à l'IFP et en particulier à l'École du Pétrole et des Moteurs, travaux de recherche ou notes de synthèse en économie, finance et gestion. La forme peut être encore provisoire, afin de susciter des échanges de points de vue sur les sujets abordés. Elle fait suite à la collection " Cahiers du CEG".

Les opinions émises dans les textes publiés dans cette collection doivent être considérées comme propres à leurs auteurs et ne reflètent pas nécessairement le point de vue de l'École ou de l'IFP.

Pour toute information complémentaire, prière de contacter :

Denis **Babusiaux** - Tél. 01 47 52 62 80



## **Résumé**

L'application des méthodes de bootstrap aux modèles de régression permet d'obtenir des approximations de la distribution des coefficients ainsi que la distribution des erreurs de prédiction. Dans cet article, nous nous intéressons à l'application des techniques de bootstrap pour déterminer des intervalles de prédiction à partir d'une modélisation économétrique où les régresseurs sont des données. Nous abordons différents problèmes liés à cette application : la détermination du nombre de répliques, le choix de la méthode de calcul de l'estimateur des moindres carrés ordinaires (pseudo-inverse ou inverse) ainsi que l'algorithme de tri de la statistique considérée. Ces investigations proviennent des besoins de prédiction des coûts dans l'industrie automobile dès la phase d'avant-projet du développement d'un nouveau véhicule. Généralement, les échantillons sont de faible taille et les termes erreur n'ont pas forcément une distribution gaussienne. Ainsi, l'utilisation des techniques de bootstrap permet d'améliorer les intervalles de prédiction en retranscrivant la distribution originale des données. Deux exemples (moteur et réservoir) illustrent la mise en œuvre de ces techniques.



## Synthèse

La réduction des coûts constitue un objectif majeur des constructeurs automobiles pour renforcer leur compétitivité sur des marchés en situation de forte concurrence. Ainsi, dans les pays industrialisés les marchés automobiles sont saturés et la compétition entre les constructeurs se fait désormais par les prix. Dans les pays émergents où les ventes sont en progression, le pouvoir d'achat des consommateurs est plus faible et les industriels doivent maîtriser leurs coûts.

L'estimation des coûts est effectuée à tous les stades de l'avant-projet à la production des véhicules. En effet, dès les premières phases de conception d'une nouvelle automobile (c'est à dire environ 36 mois avant la fabrication), les choix techniques qui sont opérés fixent une large partie des coûts futurs. Pour donner un ordre de grandeur, ce sont près de 80% des coûts qui sont fixés lorsqu'on se situe aux premiers 20% de la durée d'un projet. La détermination anticipée des coûts fait donc partie de la stratégie d'offre des constructeurs automobiles. Elle permet, de manière générale, d'aider à la conception d'une nouvelle voiture autour d'un prix cible et, en particulier, de réaliser des comparaisons entre plusieurs solutions techniques.

Les méthodes d'estimation de coût peuvent être classées en trois grandes familles : les méthodes analogiques basées sur la comparaison du nouveau projet avec des projets antérieurs, les méthodes paramétriques où le coût est relié à un ensemble de paramètres descripteurs et, finalement, les méthodes analytiques où les différents approvisionnements et temps de travail constituant le nouveau produit sont valorisés et agrégés (approche « bottom-up »). Ces dernières méthodes nécessitent une information détaillée qui n'est pas disponible au stade des avant-projets ; elles ne seront donc privilégiées que lorsque tous les éléments techniques auront été définis.

Les méthodes analogiques sont, quant à elles, utilisées dès les premières phases de développement d'un nouveau véhicule. Pour les appliquer, il faut tout d'abord définir le projet en terme de fonctions qui sont ensuite hiérarchisées<sup>(1)</sup>. Une comparaison est ensuite effectuée entre les fonctions associées au nouveau projet et celles de la (ou des) référence(s) dont les coûts sont connus. Il s'agit donc d'une approche globale qui convient bien au stade des avant-projets. Notons cependant que par comparaison il faut entendre un jugement entre les caractéristiques des deux réalisations, ce jugement devant être quantifié pour être traduisible en terme de coût. Ceci peut donc rendre délicat l'exercice d'évaluation.

Parmi les méthodes paramétriques, les formules d'estimation de coût (FEC) consistent à exprimer le coût en fonction d'un nombre restreint de paramètres techniques et à construire ainsi un modèle de régression dont on estime les coefficients. L'utilisation de la modélisation économétrique est particulièrement adaptée dans les premières étapes d'un projet automobile car elle ne nécessite pas d'information détaillée. Cependant, elle soulève des difficultés liées à la faible taille des échantillons de données et à la distribution inconnue des termes d'erreur des modèles de régression. En effet, le nombre de véhicules ou d'équipements différents proposés par un constructeur est, par nature, limité. Par ailleurs, le coût d'un équipement dépend de nombreuses considérations techniques, difficiles à préciser au stade des avant-projets, qui peuvent entraîner des surcoûts et rendre, le cas échéant, la distribution des coûts asymétrique. Dans ces conditions, les intervalles

---

<sup>1</sup> Par exemple, pour un moteur d'automobile, on peut définir la fonction « propulser l'engin », qui elle-même se décompose en plusieurs sous-fonctions « transmettre l'énergie du moteur aux roues », « transformer l'essence en énergie mécanique », etc.

standard de prédiction ne peuvent plus être utilisés ce qui enlève une grande partie à l'intérêt des résultats fournis par la modélisation.

Notre propos concerne ici la mise en œuvre du bootstrap pour établir des intervalles de prédiction sur des modèles économétriques et son application à l'estimation de modèles de coût au stade des avant-projets dans l'industrie automobile. Le principe du bootstrap consiste, en répétant un grand nombre de fois le ré échantillonnage dans les données d'origine, à construire la fonction de répartition empirique bootstrap d'une statistique considérée.

Pour un modèle économétrique, on répète ainsi un grand nombre de régressions en constituant, à chaque fois, un nouveau jeu de données par tirage avec remise dans l'échantillon de départ. On obtient ainsi une distribution empirique des coefficients et des prédictions calculées à partir de ceux-ci.

Celles-ci approchent alors de manière satisfaisante les vraies distributions des coefficients et de la prédiction qui elle, sont inconnues. Ainsi, un modèle d'estimation des coûts pourra être exploité à des fins prédictives en s'affranchissant des hypothèses inhérentes au modèle de régression linéaire.

Plusieurs aspects de la mise en œuvre des méthodes de bootstrap sont abordés de manière détaillée: ils concernent le nombre de réplifications, le calcul de l'estimateur par pseudo-inverse et le tri de la statistique considérée. Le nombre de réplifications peut être déterminé à partir des coefficients de variation de l'étendue de l'intervalle de confiance des coefficients ou de l'intervalle de prédiction lorsque ceux-ci deviennent peu variants. La mise en œuvre du bootstrap invite à privilégier le calcul des paramètres estimés en utilisant la pseudo-inverse de la matrice des variables explicatives. Cette méthode de calcul s'avère robuste en présence de multicollinéarité. Elle est également performante en terme de temps de calcul car il est suffisant de normaliser les données avant de calculer l'estimateur et non pas de les centrer et les réduire comme pour une inversion de matrice. Un algorithme modifié permet de trier rapidement la distribution empirique de la statistique considérée. Seules les queues de distribution sont triées et les éléments sont comparés aux valeurs extrêmes de celles-ci qui correspondent aux fractiles retenus.

La voiture est constituée de plusieurs milliers de pièces et accessoires, formant des sous-ensembles (ou fonctions) tels que le moteur, la caisse, la direction, etc. Pour un modèle de véhicule donné, il existe différentes versions, en termes de niveau d'équipement, de motorisation, etc. Ainsi, pour effectuer des prévisions de coûts pour tous les véhicules du modèle considéré, des modèles d'estimation de coût sont élaborés au niveau des sous-ensembles. L'entité dont le coût est modélisé correspond donc le plus souvent à un sous-ensemble de pièces assemblées.

Le périmètre des coûts étudiés concerne, ici, le coût de production, nommé Prix de Revient de Fabrication hors amortissements (PRF). Ce dernier, représentatif des dépenses engagées par l'entreprise pendant la phase de production du bien, est composé des achats (de matières et de pièces œuvrées extérieures) et d'une valeur de transformation. Le PRF représente environ 60 % du coût complet d'un véhicule (ce dernier incluant les coûts de garantie, de logistique, etc.).

La spécificité des modèles de coût implique de préciser le contexte de production industrielle. En premier lieu, les coûts des différents éléments formant les échantillons de données sont normalisés, pour un volume de fabrication moyen, représentant les conditions de production "habituelles" pour la famille de produit. Ceci permet de s'affranchir ainsi des effets d'échelle, en travaillant sur des coûts unitaires de fabrication. En second lieu, l'état de la technologie est considéré comme donné. En effet, les changements technologiques se traduisent principalement par de nouvelles machines, plus performantes ce qui ne peut pas être pris en compte dans le périmètre du coût étudié qui n'inclut pas les amortissements du capital.

Ainsi, le cadre d'analyse nous amène à spécifier des modèles de coût pour un volume de production normalisé et un état donné du système productif. Par ailleurs, nos besoins en modélisation nous conduisent à utiliser des données en "coupe transversale", c'est à dire à un instant donné puisque les phénomènes de progrès techniques ne sont pas intégrés dans la modélisation.

Deux applications permettent d'apprécier l'apport du bootstrap. Le premier exemple développé illustre l'utilisation des techniques de bootstrap sur un modèle simplifié de coût d'un moteur. L'analyse des intervalles de prédiction montre que le bootstrap permet de retranscrire l'asymétrie de la distribution des résidus dans les intervalles de prédiction. En effet, ces derniers sont décalés, par rapport aux intervalles standard, vers les plus fortes valeurs de coûts et autorisent ainsi, pour la prédiction du coût d'un nouveau moteur, des valeurs plus élevées. Ainsi, pour un moteur de 1900 cm<sup>3</sup>, l'intervalle de prédiction standard <sup>(2)</sup> est compris entre 8160 francs et 10106 francs alors que l'intervalle de prédiction bootstrap <sup>(3)</sup> est compris entre 8370 francs et 10271 francs, évitant ainsi une sous-estimation du coût de fabrication. L'utilisation des techniques de bootstrap permet donc une meilleure retranscription de l'information contenue dans l'échantillon initial, pour les intervalles de prédiction.

Le second exemple a permis d'exposer une utilisation des techniques de bootstrap adaptée, dans le cas d'une modélisation en présence de variable muette. Cette méthode permet de construire des intervalles de prédiction symétriques pour les réservoirs diesel et asymétriques vers les plus fortes valeurs de coûts pour les réservoirs essence. Ceci résulte de l'asymétrie, à la fois de la distribution de l'erreur de prédiction et de celle du surcoût lié au carburant essence. Pour des réservoirs essence de 60 litres, les intervalles de prédictions standard <sup>(2)</sup> et bootstrap <sup>(3)</sup> s'établissent respectivement, d'une part, entre 430 francs et 584 francs, et d'autre part, entre 446 francs et 619 francs. Ainsi, l'utilisation de la procédure de bootstrap des résidus stratifiés permet, lorsque l'information n'est pas disponible, de prendre en compte dans les intervalles de confiance et de prédiction, des surcoûts éventuels pour les réservoirs essence, imposés par les contraintes d'architecture du véhicule.

Les techniques de bootstrap permettent donc d'utiliser la modélisation économétrique à des fins prédictives au stade des avant-projets en retranscrivant l'asymétrie des coûts et, ceci, malgré la faible taille des échantillons de données.

---

<sup>2</sup> pour un risque de première espèce de 5 %

<sup>3</sup> pour un risque de première espèce de 5 % et en utilisant la méthode percentile-t

D'autres développements dans l'utilisation du bootstrap sont envisagés, notamment la prise en compte d'éventuelles non-linéarités dans la spécification des fonctions de coût. En effet, bien que ceci n'ait pas été mis en évidence dans les exemples traités, il faut envisager la non-linéarité de la relation comme une alternative à un modèle linéaire avec des termes d'erreur asymétriques. L'utilisation du bootstrap concerne alors les tests de spécifications ainsi que l'estimation des paramètres du modèle.



## **Introduction**

L'économétrie trouve une grande partie de ses applications industrielles dans l'obtention de prévisions. La détermination des intervalles de confiance des coefficients ainsi que des intervalles de prédiction dépend des hypothèses inhérentes aux méthodes d'estimation et, en particulier, des hypothèses sur la distribution du terme d'erreur du modèle de régression. Lorsque celles-ci ne sont plus vérifiées, les intervalles de prédiction standard ne peuvent plus être utilisés.

Le bootstrap proposé par Efron (1979) fournit une approximation d'une distribution inconnue par une distribution empirique obtenue par un processus de rééchantillonnage. L'application des méthodes de bootstrap aux modèles de régression donne ainsi une approximation de la distribution des coefficients (Freedman, 1981) et la distribution des erreurs de prédiction lorsque les régresseurs sont des données (Stine, 1985) ou des variables aléatoires (Mc Cullough, 1996).

Notre propos concerne la mise en œuvre du bootstrap pour établir des intervalles de prédiction sur des modèles économétriques lorsque les régresseurs sont connus. Ces investigations ont été suscitées par les besoins de prédiction de coûts dans l'industrie automobile dès les premières phases de développement d'un nouveau véhicule. En effet, la détermination anticipée des coûts fait partie de la stratégie d'offre des constructeurs automobiles. Elle permet, de manière générale, d'aider à la conception d'une nouvelle voiture autour d'un prix cible et, en particulier, de réaliser des comparaisons entre plusieurs solutions techniques.

Parmi les différentes approches envisageables pour effectuer des prévisions de coûts, l'utilisation de la modélisation économétrique est particulièrement adaptée dans les premières étapes d'un projet automobile car elle ne nécessite pas d'information détaillée. Cependant, elle soulève des difficultés liées à la faible taille des échantillons de données et à la distribution inconnue des termes d'erreur des modèles de régression. Dans ce contexte, le bootstrap autorise l'utilisation d'approche économétrique à des fins prédictives.

Dans la première section, nous rappelons brièvement le principe du bootstrap sur les modèles de régression. Nous abordons ensuite, dans la section 2, plusieurs problèmes liés à sa mise en œuvre : détermination du nombre de réplifications, choix de la méthode de calcul de l'estimateur des moindres carrés (pseudo-inverse ou inverse) et algorithme de tri de la statistique d'intérêt. La section suivante est consacrée à l'estimation des coûts au stade des avant-projets dans l'industrie automobile et à plusieurs applications du bootstrap (section 3). Ainsi, après l'estimation du coût d'un moteur qui ne dépend que d'une variable continue, nous présentons un modèle économétrique du coût d'un réservoir à carburant où intervient une variable binaire. Le résumé des résultats obtenus et quelques voies d'investigation forment la conclusion.

### **1 - Les techniques de bootstrap sur les modèles de régression**

Le bootstrap est une technique de rééchantillonnage basée sur des tirages aléatoires avec remise dans les données constituant un échantillon. Utilisées pour approcher la distribution inconnue d'une statistique par sa distribution empirique, les méthodes de bootstrap sont mises en œuvre afin d'améliorer la précision des estimations statistiques. Des présentations

détaillées de cette approche sont proposées, notamment par Hall (1992) et Efron et Tibshyran (1993).

L'utilisation du bootstrap sur les modèles de régression a initialement été abordé par Freedman (1981). Jeong and Maddala (1993), Vinod (1993) et Veall (1998) offrent des synthèses des nombreux développements et applications des techniques de bootstrap dans le domaine de l'économétrie qui sont ensuite apparus. Horowitz (1997) s'intéresse aux performances théoriques et numériques du bootstrap en économétrie. Nous rappelons brièvement le principe de cette méthode de rééchantillonnage ainsi que son application aux modèles de régression dans l'annexe 1.

### 1.1. Le bootstrap des résidus et les intervalles de confiance bootstrap

Le modèle de régression linéaire multiple est noté :

$$Y = X\beta + u \quad (1)$$

où  $Y$  est un vecteur  $(n,1)$ ,  $X$  une matrice  $(n,p)$ ,  $\beta$  le vecteur des coefficients à estimer  $(p,1)$  et  $u$  le vecteur des erreurs aléatoires  $(n,1)$ . Un rang d'observations  $i$  ( $i = 1, \dots, n$ ) de la matrice  $X$ , correspondant à une ligne, est noté  $X_i$   $(1,p)$ .

L'estimateur des paramètres  $\beta$  obtenu par la méthode des moindres carrés ordinaires (MCO) s'exprime comme :

$$\hat{\beta} = (X^T X)^{-1} X^T Y \quad (2)$$

et les résidus comme  $\hat{u} = Y - X\hat{\beta}$ .

Pour la suite de nos développements, nous avons retenu une approche en terme de bootstrap des résidus plutôt qu'une approche en terme de bootstrap par paire car nous ne sommes pas confrontés à un problème d'hétéroscédasticité (Flachaire, 1998).

Le modèle théorique bootstrap est le suivant :

$$Y^* = X\hat{\beta} + u^* \quad (3)$$

où  $u^*$  est un terme aléatoire issu des résidus  $\hat{u}$  de la régression initiale. A chaque itération  $b$  ( $b=1, \dots, B$ ), un échantillon  $\{y_i^*\}_{i=1}^n$ , de dimension  $(n,1)$ , est constitué à partir du modèle bootstrap (3).

Les résidus MCO étant plus petits que les erreurs qu'ils estiment, le terme aléatoire du modèle théorique bootstrap, est construit à partir des résidus transformés suivants qui sont de même norme que les termes erreurs  $u_i$  :

$$\tilde{u}_i = \frac{\hat{u}_i}{\sqrt{(1-h_i)}} - \frac{1}{n} \sum_{s=1}^n \frac{\hat{u}_s}{\sqrt{(1-h_s)}}$$

Le modèle théorique bootstrap s'exprime donc comme :

$$y_i^*(b) = X_i \hat{\beta} + \tilde{u}_i^*(b), \quad i = 1, \dots, n \quad (4)$$

Où  $\tilde{u}_i^*(b)$  est rééchantillonné à partir des  $\tilde{u}_i$ .

Soit la variable aléatoire  $z_j$  définie comme  $z_j = \frac{\hat{\beta}_j - \beta_j}{s(\hat{\beta}_j)}$ , l'intervalle de confiance standard de  $\beta_j$  découle de l'hypothèse selon laquelle  $z_j$  est distribuée selon une loi de Student à  $n-p$  degrés de liberté. Ainsi, pour un niveau de confiance  $(1 - 2\alpha)$ , cet intervalle de confiance prend alors la forme suivante :

$$\left[ \hat{\beta}_j - s(\hat{\beta}_j) \cdot t_{(1-\alpha), n-p}, \hat{\beta}_j + s(\hat{\beta}_j) \cdot t_{(\alpha), n-p} \right] \quad (5)$$

où  $t$  sont les valeurs des quantiles  $(\alpha)$  et  $(1 - \alpha)$  de la distribution de Student à  $n-p$  degrés de liberté.

Les intervalles de confiance bootstrap sont construits à partir des deux approches percentile et percentile-t. La première méthode, basée uniquement sur les estimations bootstrap, est la méthode la plus simple d'obtention d'intervalles de confiance. Pour un niveau  $(1 - 2\alpha)$ , l'intervalle de confiance percentile pour le paramètre  $\beta_j$  est donné par :

$$\left[ \hat{\beta}_j^*(\alpha B), \hat{\beta}_j^*((1 - \alpha)B) \right] \quad (6)$$

où  $\hat{\beta}_j^*(\alpha B)$  représente la  $\alpha B$ -ième valeur (respectivement  $\hat{\beta}_j^*((1 - \alpha)B)$  la  $(1 - \alpha)B$ -ième valeur) de la liste ordonnée des  $B$  réplifications bootstrap. Les valeurs seuils sont donc choisies telles que  $\alpha$  % des réplifications ont fourni des  $\hat{\beta}_j^*$  plus petits (grands) que la borne inférieure (supérieure) de l'intervalle de confiance percentile.

La procédure bootstrap percentile-t consiste à estimer la fonction de répartition de  $z_j$  directement à partir des données. Cela revient à construire une table statistique à partir de la fonction de répartition empirique des  $B$  réplifications bootstrap  $z_j^*$ . Cette table est nommée table bootstrap. Les  $z_j^*$  sont définies comme :

$$z_j^* = \frac{\hat{\beta}_j^* - \hat{\beta}_j}{s^*(\hat{\beta}_j^*)} \quad (7)$$

Soit  $\hat{F}_{z_j^*}$  la fonction de répartition empirique des  $z_j^*$ , le fractile à  $\alpha$  %,  $\hat{F}_{z_j^*}^{-1}(\alpha)$ , est estimé par la valeur  $\hat{t}^{(\alpha)}$  telle que  $\# \{z_j^*(b) \leq \hat{t}^{(\alpha)}\} / B = \alpha$ .

Finalement, l'intervalle de confiance percentile-t pour  $\beta_j$  s'écrit :

$$\left[ \hat{\beta}_j - s(\hat{\beta}_j) \cdot \hat{t}^{(1-\alpha)}, \hat{\beta}_j + s(\hat{\beta}_j) \cdot \hat{t}^{(\alpha)} \right] \quad (8)$$

Ainsi, l'intervalle de confiance percentile-t est l'analogue bootstrap de l'intervalle de confiance standard.

## 1.2. Les intervalles de prédiction bootstrap

A la suite de Stine (1985) et Breiman (1992), notre cadre de travail nous conduit à utiliser le bootstrap pour la construction des intervalles de prédiction sur des modèles de régression avec des régresseurs fixes, dont les valeurs sont connues (prévision non conditionnelle). Notons cependant que la construction des intervalles de prédiction bootstrap sur des modèles avec des régresseurs stochastiques est proposée par McCullough (1996).

Pour un nouveau rang  $f$  d'observations des variables explicatives  $X_f$ , la prédiction de coût  $\hat{y}_f$  est calculée à partir du modèle de régression :  $\hat{y}_f = X_f \hat{\beta}$ .

L'intervalle de prédiction standard découle, comme les intervalles de confiance des coefficients de la régression, de l'hypothèse de Normalité des erreurs. Ainsi, pour un niveau de confiance  $(1 - 2\alpha)$ , cet intervalle de prédiction standard s'écrit :

$$\left[ \hat{y}_f - s_f \cdot t_{(1-\alpha), n-p}, \hat{y}_f + s_f \cdot t_{(\alpha), n-p} \right] \quad (9)$$

L'utilisation du bootstrap, pour préciser les intervalles de prédiction, conduit à étudier la distribution de l'erreur de prédiction. Aussi, afin de conserver le même processus générateur de données (PGD) pour les estimations des coefficients et des prédictions, les intervalles de prédiction bootstrap sont obtenus avec la procédure du bootstrap des résidus. De manière similaire à la construction des intervalles de confiance, il existe deux principales méthodes de construction des intervalles de prédiction bootstrap : l'approche percentile et percentile-t.

– L'intervalle de prédiction percentile

La méthode percentile consiste à utiliser l'approximation bootstrap de la distribution de l'erreur de prédiction :  $e_f = \hat{y}_f - y_f$ , pour construire un intervalle de prédiction de  $y_f$ .

Les répliques bootstrap de la future valeur  $y_f^*$ , pour le nouveau rang d'observations  $X_f$  sont générées suivant le même modèle (4) :

$$y_f^* = X_f \hat{\beta} + \tilde{u}_f^* \quad (10)$$

Le terme d'erreur  $\tilde{u}_f^*$  est issu, comme les  $\tilde{u}^*$ , d'un tirage avec remise dans la distribution empirique des résidus transformés.

Pour chacune des  $B$  répliques bootstrap, nous calculons l'estimateur bootstrap. Ainsi, la prévision et l'erreur de prédiction bootstrap s'écrivent respectivement :

$$\begin{aligned}\hat{y}_f^*(b) &= X_f \hat{\beta}^*(b) \\ e_f^*(b) &= \hat{y}_f^*(b) - y_f^*(b)\end{aligned}\tag{11}$$

En utilisant l'équation (9), nous pouvons réécrire l'erreur de prédiction bootstrap comme :

$$e_f^* = \hat{y}_f^* - \hat{y}_f - \tilde{u}_f^*\tag{12}$$

Cette dernière dépend donc, par nature, de la prédiction MCO initiale  $\hat{y}_f$ .

Les  $B$  répliques bootstrap de l'erreur de prédiction fournissent la distribution empirique de  $e_f^* : G^*$ . Les quantiles de cette distribution empirique, notés  $G^{*-1}(1-\alpha)$  et  $G^{*-1}(\alpha)$ , sont alors utilisés pour construire un intervalle de prédiction bootstrap.

Un intervalle de prédiction percentile est finalement de la forme suivante :

$$\left[ \hat{y}_f - G^{*-1}(1-\alpha); \hat{y}_f - G^{*-1}(\alpha) \right]\tag{13}$$

– L'intervalle de prédiction percentile-t

De manière identique à l'intervalle de confiance, la construction de l'intervalle de prédiction, avec la méthode percentile-t implique le calcul, pour chaque échantillon bootstrap, de l'estimateur bootstrap de l'écart-type. Ainsi, pour établir des intervalles de prédiction percentile-t, l'estimateur bootstrap de l'écart-type de prédiction est nécessaire, pour chacune des répliques. Ce dernier s'écrit :

$$s_f^* = s^* \cdot \sqrt{(1+h_f)}\tag{14}$$

où  $s^*$  est l'estimateur bootstrap de l'écart-type des termes erreurs et  $h_f = X_f (X^T X)^{-1} X_f^T$ .

La procédure percentile-t consiste à construire les statistiques  $z_f^*$ , telles que :

$$z_f^* = \frac{e_f^*}{s_f^*} = \frac{\hat{y}_f^* - \hat{y}_f - \tilde{u}_f^*}{s_f^*}\tag{15}$$

La distribution bootstrap de  $z_f^*$  définit l'intervalle de prédiction bootstrap percentile-t. Les quantiles  $z_{f(1-\alpha)}^*$  et  $z_{f(\alpha)}^*$ , remplacent ainsi les valeurs critiques de la distribution de Student, prises en compte dans l'intervalle de prédiction standard (cf.. équation 9).

Un intervalle de prédiction percentile-t s'écrit donc :

$$\left[ \hat{y}_f - s_f \cdot z_{f(1-\alpha)}^*; \hat{y}_f - s_f \cdot z_{f(\alpha)}^* \right] \quad (16)$$

Notons que, comme pour l'intervalle de confiance des coefficients, le quantile  $(1-\alpha)$  de la distribution de  $z_f^*$  définit la borne inférieure de l'intervalle de prédiction et inversement pour le quantile  $(\alpha)$ .

Une distribution symétrique de  $z_f^*$  implique donc la symétrie de l'intervalle de prédiction percentile-t. Cependant, dans le cas contraire, l'asymétrie est retranscrite de manière inversée, pour ce dernier. Par exemple, si  $z_f^*$  possède une queue de distribution plus longue vers la droite, les quantiles  $z_{f(1-\alpha)}^*$  et  $z_{f(\alpha)}^*$  sont décalés vers les valeurs élevées des erreurs de prédiction bootstrap, comparativement aux quantiles correspondants d'une distribution symétrique. L'intervalle de prédiction percentile-t résultant est donc décalé vers la gauche, asymétrique autour de la valeur prédite MCO. Ainsi, sa construction implique une sorte de "correction automatique du biais" et permet l'acceptation, pour un niveau de confiance donné, de valeurs prédites plus faibles que l'intervalle de prédiction standard, symétrique.

## 2 - La mise en œuvre des méthodes

Les développements que nous proposons pour la mise en œuvre du bootstrap s'inscrivent dans le prolongement de Efron et Tibshirani (1993), Booth et Sarkar (1998), Davidson et McKinnon (1998) pour la détermination du nombre de réplifications bootstrap. À la suite de Mc Cullough et Vinod (1996), nous nous intéressons au choix de la méthode de calcul de l'estimateur MCO par pseudo-inverse ainsi qu'à l'algorithme de tri de la statistique d'intérêt.

### 2.1. Le nombre de réplifications bootstrap B

Efron and Tibshyran (1993) préconisent d'effectuer un nombre de réplifications bootstrap de l'ordre de 25 pour le calcul de l'écart-type d'un estimateur et de l'ordre du millier pour les intervalles de confiance bootstrap.

La construction de l'intervalle de confiance (IC)<sup>4</sup> bootstrap implique les fractiles (2,5 %) et (97,5 %) de la distribution empirique bootstrap de la statistique  $z^*$ . Nous nous sommes intéressé à l'incidence du nombre de réplifications sur la détermination des fractiles afin de définir le nombre minimum de réplifications nécessaires avant d'obtenir des fractiles qui varient peu.

---

<sup>4</sup> Le processus, présenté pour la construction de l'IC bootstrap, s'applique de manière identique pour la construction de l'intervalle de prédiction bootstrap.

Intuitivement, il paraît logique que l'estimation des fractiles d'une distribution nécessite un nombre plus élevé d'échantillons bootstrap que le calcul de l'écart-type de l'estimateur, par exemple. En effet, cette dernière dépend de la queue de distribution, où peu d'échantillons apparaissent. La question revient alors à déterminer le nombre de réplifications bootstrap à partir duquel la valeur du fractile peut être considérée comme stable.

Le processus mis en oeuvre consiste, pour un certain nombre de valeurs de  $B$ , à effectuer un ensemble de simulations (noté  $k$ ), visant à juger de la stabilité des résultats, lorsque les racines (valeurs de départ) du générateur de nombre pseudo-aléatoires sont différentes.  $k = 100$  simulations ont été réalisées, ce qui semble raisonnable, compte tenu de nos investigations. Par ailleurs, les simulations ont été effectuées pour les nombres de réplifications bootstrap suivants :  $B = 20, 30, 100, 500, 1\ 000, 5\ 000$  et  $10\ 000$ .

Plutôt que d'étudier la variabilité de chacun des fractiles (2,5 %) et (97,5 %) de la distribution séparément, nous considérons l'étendue de l'intervalle entre ces derniers, que nous nommons "intervalle de confiance de  $z^*$ ". Ainsi, nos travaux s'appuient sur l'analyse de la variabilité des étendues de ces IC, sur 100 simulations, en fonction de  $B$ .

En résumé, pour chacune des  $k$  simulations,  $k = 1, \dots, 100$ , les  $B$  réplifications bootstrap fournissent la distribution empirique de  $z^*$ , à partir de laquelle sont extraits les fractiles d'intérêt. Pour un nombre de réplifications  $B$  donné, les simulations permettent donc d'obtenir 100 IC de  $z^*$  et leurs étendues. Leur variabilité est ensuite étudiée.

Les comparaisons des distributions des étendues sont effectuées deux à deux, pour des valeurs croissantes de  $B$ . Pour ce faire, les caractéristiques de valeurs centrales (moyenne, médiane) et de dispersion sont calculées et trois critères sont examinés : l'égalité des médianes<sup>5</sup>, l'égalité des variances et le coefficient de variation (CV).

Les deux premiers critères font l'objet de tests statistiques, respectivement le test non paramétrique de Wilcoxon d'égalité des médianes et le test de Fisher d'égalité des variances de deux échantillons indépendants. L'évolution du CV en fonction du nombre de réplifications fait l'objet de l'analyse du dernier critère. L'application et l'interprétation de cette procédure sont présentées sur l'exemple du modèle d'estimation de coût des moteurs.

## 2.2. Le calcul de l'estimateur des paramètres de la régression

Le calcul de l'estimateur MCO est effectué, classiquement, à partir de la formule (2) en inversant la matrice  $(X^T X)$ . L'inversion de cette matrice soulève des problèmes d'instabilité numérique lorsque la matrice  $X$  des variables explicatives est mal conditionnée (Belsley, Kuh et Welsch, 1980). Il est préférable de calculer l'estimateur MCO à partir de la décomposition en valeurs singulières :

$$\begin{matrix} X & = & U & D & V^T \\ (n,p) & & (n,p) & (p,p) & (p,p) \end{matrix} \quad (17)$$

---

<sup>5</sup> Nous avons retenu la médiane comme indicateur de valeur centrale car, contrairement à la moyenne, il est insensible aux variations des valeurs extrêmes de la distribution.

D est la matrice diagonale des valeurs singulières de X. U est la matrice orthogonale des p vecteurs propres associés aux p valeurs propres non nulles de  $(XX^T)$  et V est la matrice orthogonale des vecteurs propres de  $(X^T X)$ . En notant  $X^+ = V D^+ U^T$  la pseudo-inverse de X, l'estimateur MCO s'exprime comme :

$$\hat{\beta} = X^+ y \quad (18)$$

et sa matrice de variance-covariance estimée s'écrit :

$$\hat{V}(\hat{\beta}) = s^2 V D^{-2} V^T$$

La matrice de projection qui lie les résidus et le terme d'erreur s'exprime simplement comme :

$$I_n - X(X^T X)^-1 X^T = I_n - U U^T \quad (19)$$

Nous avons comparé les performances obtenues en calculant classiquement l'estimateur MCO par inversion de matrice (2) et par pseudo-inverse (18). Après un ensemble de tests sur Matlab (cf. Juan, 1999), les deux méthodes ont été programmées directement en Fortran. L'algorithme le mieux adapté pour l'inversion de la matrice  $(X^T X)$  qui est définie positive repose sur une décomposition de Cholesky (Seak, 1972). Une analyse de ses performances numériques est fournie dans Lantz (1983). Le calcul de la pseudo-inverse  $X^+$  a été réalisé en reprenant l'algorithme de Golub et Reinsch qui est décrit dans Forsythe et al. (1977).

La multicollinéarité a moins de conséquences sur le conditionnement de la matrice X lorsque l'on considère des données normées ou centrées et réduites (Belsley et al., 1980 – Belsley, 1984 – Erkel-Rousse, 1995). Nous avons testé l'incidence de ces transformations pour différentes dimensions de la matrice des variables explicatives ( $n=10, \dots, 50$ ,  $p=2, \dots, 10$ ). Pour chaque dimension nous avons généré 1000 matrices X telles que  $X_1^T$  suit une loi uniforme (0,1) et les vecteurs j suivants sont construits comme  $X_j^T = X_{j-1}^T + v$  où v suit une loi uniforme (0, 0.001). Nous avons mesuré les erreurs de calcul comme la somme des valeurs absolues des écarts entre  $XX^+$  ou  $(X^T X)(X^T X)^{-1}$  et la matrice identité. Le tableau 1 résume ces écarts pour différentes dimension (n,p)<sup>6</sup>. Ils conduisent à estimer les vecteur des paramètres par pseudo-inverse sur une matrice X normée à 1 ou le cas échéant à le calculer par inversion de matrice sur des données centrées et réduites lorsque p est faible.

Tableau 1 – Erreur de calcul sur les produits de matrice  $XX^+$  et  $(X^T X)(X^T X)^{-1}$

Dimension de X	$XX^+$ données normées à 1	$(X^T X)(X^T X)^{-1}$ données normées à 1	$(X^T X)(X^T X)^{-1}$ données centrées et réduites
(50,2)	$<10^{-10}$	$<10^{-10}$	$<10^{-10}$
(50,4)	$<10^{-10}$	$0.61 \times 10^{-6}$	$0.11 \times 10^{-6}$
(50,8)	$<10^{-10}$	$0.54 \times 10^{-5}$	$0.59 \times 10^{-6}$

<sup>6</sup> L'ensemble des résultats est disponible auprès des auteurs.



(50,10)	$<10^{-10}$	$0.50 \times 10^{-4}$	$0.50 \times 10^{-5}$
---------	-------------	-----------------------	-----------------------

Nous avons ensuite comparé les temps de calcul pour ces deux méthodes sur le bootstrap des résidus (en utilisant l'algorithme de tri décrit dans le paragraphe suivant). Nous avons généré 1000 échantillons de données pour différentes dimensions (n,p) comme précédemment. Pour chaque échantillon, nous avons appliqué un bootstrap des résidus et la méthode percentile avec 1000 réplifications. Le tableau 2 donne les temps de calcul obtenus sur un micro-ordinateur de type PC (Pentium II, 300 Mhz). L'utilisation d'une pseudo-inverse divise les temps de calcul par un facteur proche de deux et évitant le centrage des données.

Tableau 2 – Temps de calcul pour le bootstrap des résidus et la méthode percentile sur 1000 réplifications

Dimension de X	$X^+$ données normées à 1	$(X^T X)^{-1}$ données centrées et réduites
(10,2)	0.0103	0.0208
(50,4)	0.0518	0.0972
(90,6)	0.1001	0.1823

Unité : seconde

### 2.3 L'algorithme de tri des réplifications

Les deux algorithmes de tri les plus répandus sont le tri par comparaison de tous les éléments et le tri par partition. Cette dernière méthode est privilégiée pour de grands échantillons puisqu'elle requiert de l'ordre de  $B \log_2 B$  opérations pour trier les B réplifications alors que la méthode classique nécessite de l'ordre de  $B^2/2$  opérations.

Le tri modifié que nous proposons consiste à ne s'intéresser qu'à la détermination des fractiles  $\alpha$  et  $1-\alpha$ . Ainsi, dans la phase de tri, nous traitons uniquement  $(\alpha B + 1)$  réplifications bootstrap, à la différence de l'algorithme classique, dans lequel la totalité des réplifications (B) est triée. Il peut être décomposé en deux étapes.

La première consiste à construire un vecteur  $S_1$  de dimension  $(\alpha B + 1)$ , constitué des  $(\alpha B + 1)$  premières statistiques  $\hat{\theta}^*(b)^7$ , classées par ordre croissant. Ainsi,  $S_1(1)$  et  $S_1(\alpha B + 1)$  correspondent respectivement à la plus petite et la plus grande valeur des  $(\alpha B + 1)$  premières statistiques bootstrap.  $S_1$  est ensuite dupliqué dans un second vecteur  $S_2$ , de dimension identique.

Dans une deuxième étape, chacune des  $B - (\alpha B + 1)$  statistiques suivantes est comparée avec les éléments de  $S_1$ . Le vecteur  $S_1$  est utilisé pour la recherche des plus petits éléments, parmi les B réplifications bootstrap, suivant la procédure suivante. Cette dernière est appliquée pour chaque itération bootstrap  $b$ ,  $b = \alpha B + 1, \dots, B$ .  $\hat{\theta}^*(b)$  est tout d'abord comparé au plus grand élément de  $S_1$ . S'il est plus grand ou égal à celui-ci, on passe à la comparaison avec  $S_2$ . Dans le cas contraire,  $\hat{\theta}^*(b)$  est alors comparé à chacun des éléments du vecteur  $S_1$ . S'il est inférieur ou égal à l'élément k de  $S_1$ ,  $S_1(k)$  est remplacé par

<sup>7</sup> Elles correspondent aux valeurs de la statistique bootstrap, pour les  $b=1, \dots, (\alpha B + 1)$  premières itérations.

$\hat{\theta}^*(b)$  et pour chaque indice supérieur à  $k$  de ce vecteur, les éléments sont décalés d'un rang, vers le rang supérieur. Finalement, le vecteur  $S_1$  contient les  $(\alpha B + 1)$  plus petites valeurs des  $B$  statistiques  $\hat{\theta}^*(b)$ . Ainsi, le quantile  $\alpha$  correspond à l'élément  $S_1(\alpha B)$

De manière similaire on compare les  $B - (\alpha B + 1)$  statistiques bootstrap avec les éléments de  $S_2$ . Il s'agit donc du vecteur utilisé pour la recherche des plus grands éléments, parmi les  $B$  répliquions bootstrap. Ainsi, le quantile  $(1 - \alpha)$  de la distribution bootstrap correspond à l'élément  $S_2(1)$ . Notons qu'il est nécessaire de classer  $(\alpha B + 1)$  répliquions bootstrap et non pas  $(\alpha B)$ , pour extraire le quantile  $(1 - \alpha)$ . En effet, ce dernier correspond au premier élément du vecteur  $S_2$ , de dimension  $(\alpha B + 1)$ .

Dans cette version améliorée de l'algorithme, deux étapes sont distinguées : le tri des  $(\alpha B + 1)$  premières répliquions bootstrap puis, pour chaque répliquion suivante, les comparaisons avec les éléments de  $S_1$  et  $S_2$ . Le nombre de comparaisons de la première étape est égal à  $\frac{\alpha B(\alpha B + 1)}{2}$ . Celui de la seconde étape peut être encadré par un cas minimum, où chaque répliquion est toujours inférieure à  $S_2(1)$  et supérieure à  $S_1(\alpha B + 1)$  et un cas maximum où elle est toujours comparée aux  $(\alpha B + 1)$  éléments de  $S_1$  et de  $S_2$ . Il est donc compris entre  $[B - (\alpha B + 1)] \times 2$  et  $[B - (\alpha B + 1)] \times (\alpha B + 1) \times 2$ .

Le nombre total de comparaisons, dans la version améliorée de l'algorithme de tri, est compris entre  $[(\alpha B + 1) \times (\alpha B - 1)] + B$  et  $(\alpha B + 1) \times (B - 1)$ , il est d'ordre  $B^2$ . Asymptotiquement, le nombre d'opérations de cet algorithme est du même ordre qu'un tri classique. Cependant, pour  $B = 5\,000$  et  $\alpha = 0,025$ , le nombre d'opérations effectuées est de l'ordre de 12 millions avec l'algorithme classique et compris entre 20 000 et 600 000, avec la version améliorée de l'algorithme. Ainsi, cette dernière permet de diviser le nombre de comparaisons au moins par 20.

### 3 - Applications à l'estimation des coûts dans le secteur automobile

Le marché automobile est dans une situation de forte concurrence. La réduction des coûts constitue un objectif majeur des firmes pour renforcer leur compétitivité. Ainsi, dans les pays industrialisés les marchés automobiles sont saturés et la compétition entre les constructeurs se fait désormais par les prix. F. Verboven (1996) propose une analyse du marché automobile européen en terme de concurrence oligopolistique. Dans les pays émergents où les ventes sont en progression, le pouvoir d'achat des consommateurs est plus faible et les industriels doivent maîtriser leurs coûts.

L'estimation des coûts est effectuée à tous les stades de l'avant-projet à la production des véhicules. En effet, dès les premières phases de conception d'une nouvelle automobile (c'est-à-dire environ 36 mois avant la fabrication), les choix techniques qui sont opérés fixent une large partie des coûts futurs (Juan, 1999). La prédiction des coûts revêt donc un enjeu majeur.

L'économétrie fournit une méthode d'estimation et de prédiction des coûts particulièrement utile à cette fin. En effet, les modèles économétriques permettent d'expliquer les coûts en

fonction des quelques paramètres techniques qui constituent la seule information disponible au démarrage d'un projet.

Cependant l'estimation de tels modèles soulève des difficultés liées à la faible taille des échantillons de données et de la distribution asymétrique des erreurs dans les modèles de régression. Les méthodes de bootstrap permettent de pallier ces difficultés en fournissant une approximation de la distribution des erreurs de prédiction par leur distribution empirique. Nous illustrons ceci au travers de deux exemples. Au préalable, nous définissons le coût qui est étudié et nous présentons les données qui sont utilisées.

### **3.1. Les données et la forme générale des modèles**

La voiture est constituée de plusieurs milliers de pièces et accessoires, formant des sous-ensembles (ou fonctions) tels que le moteur, la caisse, la direction, etc. Pour un modèle de véhicule donné, il existe différentes versions, en termes de niveau d'équipement, de motorisation, etc. Ainsi, pour effectuer des prévisions de coûts pour tous les véhicules du modèle considéré, des modèles d'estimation de coût sont élaborés au niveau des sous-ensembles. L'entité dont le coût est modélisé correspond donc le plus souvent à un sous-ensemble de pièces assemblées.

Par ailleurs, il importe de préciser le périmètre des coûts étudiés. Notre travail porte sur le coût de production, nommé Prix de Revient de Fabrication hors amortissements (PRF). Ce dernier, représentatif des dépenses engagées par l'entreprise pendant la phase de production du bien, est composé des achats (de matières et de pièces oeuvrées extérieures) et d'une valeur de transformation. Le PRF représente environ 60 % du coût complet d'un véhicule (ce dernier incluant les coûts de garantie, de logistique, etc.). Notons que le PRF hors amortissements n'inclut pas les amortissements du capital.

La spécificité des modèles de coût implique de préciser le contexte de production industrielle. En premier lieu, les coûts des différents éléments de la base de données sont normalisés, pour un volume de fabrication moyen, représentant les conditions de production "habituelles" pour la famille de produit. Nous nous affranchissons ainsi des effets d'échelle en travaillant sur des coûts unitaires de fabrication.

En second lieu, l'état de la technologie est considéré comme donné. En effet, les changements technologiques se traduisent principalement par de nouvelles machines, plus performantes. Or, le périmètre du coût étudié n'inclut pas les amortissements du capital.

Ainsi, le cadre d'analyse nous amène à spécifier des modèles de coût pour un volume de production normalisé et un état donné du système productif. Par ailleurs, nos besoins en modélisation nous conduisent à utiliser des données en "coupe transversale" du coût de chacun des éléments de la base de données à l'instant  $t$ . Les estimations sont donc effectuées dans un cadre statique et n'intègrent pas les phénomènes de progrès techniques.

Le coût est donc expliqué par les caractéristiques techniques pertinentes du produit, à un instant donné de la technologie, à des fins de prédiction de coût pour de nouveaux produits. Notre démarche est donc différente de celle utilisée dans l'application des techniques de bootstrap sur les modèles de frontières de production (Simar, 1992). En effet, dans ce cas il s'agit d'évaluer l'efficacité d'une unité de production par rapport à une frontière efficace, à

l'aide du bootstrap. Finalement, notons que si aucun *a priori* n'est posé, quant à la forme de la relation liant le coût aux descripteurs, les formes linéaires ou multiplicatives sont le plus souvent retenues

### 3.2. Le modèle d'estimation de coût des moteurs

L'exemple porte sur le périmètre du moteur, avec ses composants électriques. De manière très schématique, le moteur est composé du carter cylindre, de l'attelage mobile (vilebrequin, volant moteur, etc.), de pièces assurant la distribution (arbre à cames, soupapes, etc.), et de la culasse (y compris son couvercle). Les équipements électriques comprennent le démarreur, la bobine et les bougies d'allumage, l'alternateur, etc.

Le moteur constitue, avec l'assemblage final du véhicule, une fonction dont la fabrication est, le plus souvent, interne à l'entreprise. Les données de coûts sont donc des PRF hors amortissements. Par ailleurs, le moteur représente une partie non négligeable (de l'ordre de 20 %) du PRF d'un véhicule.

Il existe différentes familles de moteurs, distinguées principalement par le mode de carburation (essence ou diesel), le mode d'injection (directe, monopoints, etc.) et la cylindrée. L'étude porte sur des moteurs de cylindrée supérieure à 1 700 cm<sup>3</sup>, ces derniers formant un échantillon de travail homogène de 15 moteurs.

Notons, dans le périmètre technique du moteur, la présence des composants supplémentaires (support, tuyaux, etc.), imposés par des équipements tels que la direction assistée ou le conditionnement d'air. Ceux-ci sont présents ou non sur le moteur, en fonction des prestations à assurer par le véhicule sur lequel ils sont destinés à être montés.

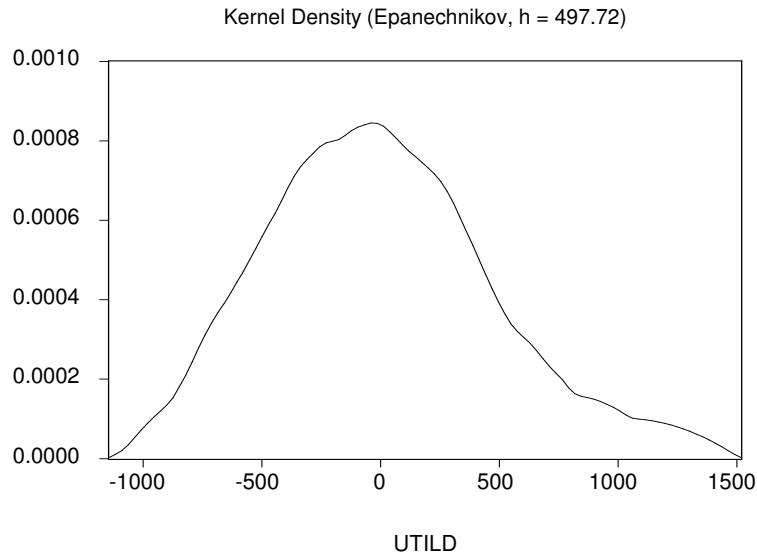
Dans l'équation économétrique le coût d'un moteur est exprimé comme une fonction de sa cylindrée. Le modèle possède une variable explicative et un terme constant ; il est estimé par la méthode des moindres carrés ordinaires (MCO). Le R<sup>2</sup> est égal à 0.97 et les coefficients estimés<sup>8</sup> de la régression sont significativement différents de zéro, pour un risque de première espèce de 5 %. Le test de White F(2,12) = 0,66 ne permet pas de rejeter l'hypothèse d'homoscédasticité.

Les résidus transformés de la régression  $\tilde{u}$ , à partir desquels sont constitués les échantillons dans la procédure de bootstrap des résidus, sont illustrés dans la figure 1. Nous remarquons une queue de distribution plus longue vers la droite, traduisant la présence de résidus positifs élevés. En effet, le résidu "extrême" correspondant à un moteur de 1 783 cm<sup>3</sup>, que nous notons  $\hat{u}_6$ , possède la valeur la plus élevée de l'ensemble des résidus. Ce moteur, ainsi qu'un second dans l'échantillon, sont équipés du conditionnement d'air et correspondent aux résidus les plus élevés de la régression

---

<sup>8</sup> Pour des raisons de confidentialité, les valeurs des paramètres estimés ne sont pas reproduites.

Figure1 L'estimateur à noyau de la densité des résidus transformés



La procédure de détermination du nombre de réplifications bootstrap nécessaires est ensuite mise en oeuvre (cf. paragraphe 2.1). Les résultats présentés portent sur l'exemple du coefficient de la cylindrée  $z_1^*$ .

– Le nombre de réplifications  $B$

Tableau 3 Les tests des étendues des IC de  $z_1^*$

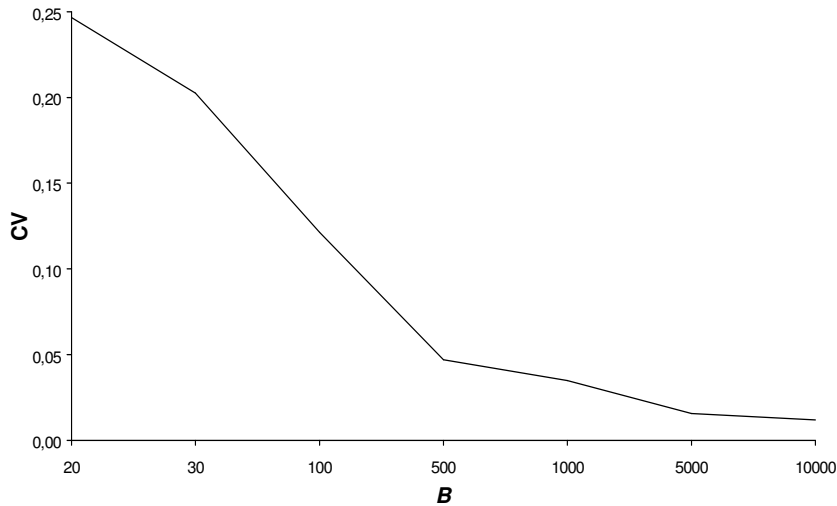
<b>B</b>	<b>Médiane des étendues</b>	<b>Test de Wilcoxon</b>	<b>P-marg.</b>	<b>Écart-type des étendues</b>	<b>Test de Fisher</b>	<b>P-marg.</b>
<b>20</b>	3,910			0,973		
<b>30</b>	4,307	3,676	0,000	0,889	1,197	0,372
<b>100</b>	4,224	0,864	0,388	0,520	2,924	0,000
<b>500</b>	4,273	1,130	0,258	0,203	6,565	0,000
<b>1 000</b>	4,304	0,458	0,647	0,150	1,830	0,003
<b>5 000</b>	4,301	0,030	0,976	0,067	5,019	0,000
<b>10 000</b>	4,317	1,604	0,109	0,051	1,727	0,007

Les résultats des tests de Wilcoxon pour l'égalité des médianes et des tests de Fisher pour l'égalité des variances des étendues des intervalles de confiance de  $z_1^*$  ne permettent pas d'établir de manière certaine une valeur de  $B$  à partir de laquelle nous pourrions accepter l'hypothèse d'égalité des médianes et des variances des IC (cf. tableau 3).

L'évolution du CV en fonction du nombre de réplifications est représentée sur la figure 2. La décroissance de CV est forte, pour les faibles valeurs de  $B$ , puis s'affaiblit progressivement. Ainsi, à partir de  $B = 5\ 000$ , il se stabilise. Ce troisième critère nous permet finalement de sélectionner le nombre de réplifications bootstrap  $B = 5\ 000$ , à partir duquel les étendues des IC de  $z_1^*$  (et donc les fractiles (2,5 %) et (97,5 %) de la distribution) peuvent être

considérées comme stables. De manière similaire,  $B = 5\,000$  répliques sont retenues pour la construction des intervalles de prédiction bootstrap.

Figure 2 Le coefficient de variation des étendues des IC de  $z_1^*$



– La prédiction de coût

L'ajustement de la régression est maintenant utilisé pour prévoir les PRF d'un nouveau moteur, de cylindrée égale à  $1\,900\text{ cm}^3$ . Nous reportons, pour les intervalles de prédiction bootstrap, les deux méthodes de construction : percentile et percentile-t (cf. tableau 4). Les intervalles percentile possèdent une étendue nettement plus réduite que les intervalles percentile-t ou standard. Ainsi, la méthode percentile conduit à des intervalles trop "optimistes" (trop petits) et n'apparaît pas pertinente pour ce type d'investigation. En effet, l'erreur de prédiction n'est pas une statistique pivot et l'inférence bootstrap peut s'avérer erronée dans ce cas. Notons que les intervalles de prédiction bootstrap sont décalés vers les plus fortes valeurs de coût (forme supérieure à 1) par rapport aux intervalles standard.

Tableau 4 Les prévisions MCO et les intervalles de prédiction standard et bootstrap

Moteur Cyl. en $\text{cm}^3$	Prévision MCO	Intervalle de prédiction standard			
		2,5 %	97,5 %	Étendue	Forme <sup>9</sup>
1900	9133,42	8160,67	10106,20	1945,53	1,00

Cyl. en $\text{cm}^3$	Intervalle de prédiction percentile				Intervalle de prédiction percentile			
	2,5 %	97,5 %	Étendue	Forme	2,5 %	97,5 %	Étendue	Forme
1900	8419,66	10179,57	1759,91	1,47	8370,00	10271,97	1901,96	1,49

Unités : francs français

<sup>9</sup> La "forme" est définie comme :  $\frac{\text{sup} - \hat{y}_f}{\hat{y}_f - \text{inf}}$ , avec inf et sup correspondant respectivement aux bornes inférieures et supérieures de l'intervalle de prédiction bootstrap,  $\hat{y}_f$  étant la prévision MCO.

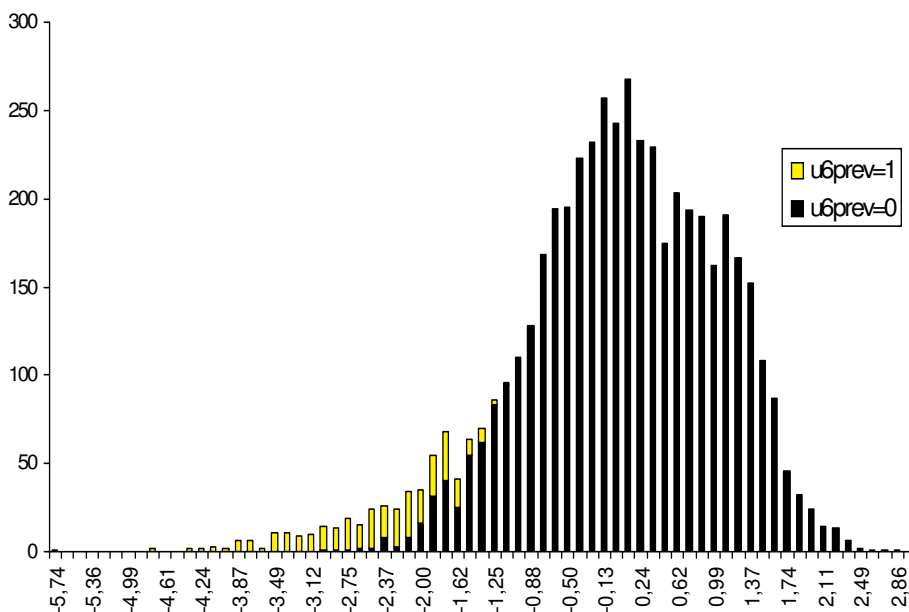
– Le repérage du résidu "extrême" dans les répliques bootstrap pour la prédiction

Afin d'expliquer l'asymétrie des intervalles de prédiction bootstrap, nous étudions l'impact du tirage du résidu extrême  $\hat{u}_6$  dans le PGD bootstrap, sur la distribution de l'erreur de prédiction bootstrap. Pour ce faire, nous examinons, pour chacune des répliques, le résidu  $\tilde{u}_f^*$  du modèle théorique bootstrap de la prédiction et vérifions s'il correspond ou pas au résidu  $\hat{u}_6$ . Le tableau 5 présente, pour ces deux possibilités, les valeurs moyennes et écart-types de  $z_f^*$  (la statistique bootstrap de l'erreur de prédiction normée), pour les cinq mille répliques. La figure 3 illustre, pour la prédiction du coût d'un moteur de 1900 cm<sup>3</sup>, la distribution de l'erreur de prédiction bootstrap, en distinguant les cas où le résidu extrême est tiré.

Tableau 5 - Les caractéristiques de  $z_f^*$

Moteur	$z_f^*$	$\tilde{u}_f^* \neq \hat{u}_6$	$\tilde{u}_f^* = \hat{u}_6$	Total
1900 cm <sup>3</sup>	Moyenne	0,144	-2,472	-0,022
	Écart-type	0,896	0,745	1,092
	Effectifs	4682	318	5000

Figure 3 La distribution de la statistique  $z_f^*$ , pour un moteur de 1 900 cm<sup>3</sup>



La distribution bootstrap de l'erreur de prédiction normée  $z_f^*$  paraît fortement asymétrique vers les valeurs négatives. De plus, nous constatons (figure 3) que cette asymétrie est due au tirage du résidu extrême  $\hat{u}_6$ , dans le modèle théorique de prévision bootstrap. En effet, comme  $e_f^* = \hat{y}_f^* - \hat{y}_f - \tilde{u}_f^*$ , si  $\tilde{u}_f^* = \hat{u}_6$  (forte valeur positive), l'erreur de prédiction bootstrap est fortement négative, ce que nous retrouvons dans l'histogramme. Ainsi, l'intervalle de prédiction résultant est décalé vers les plus grandes valeurs de coût, puisque



le fractile (2,5 %) de la distribution détermine la borne supérieure de l'intervalle et vice versa pour le fractile (97,5 %). Notons que ce phénomène se répète pour la distribution de la statistique  $z_f^*$  effectuée pour d'autres prédictions (cf. Juan, 1999).

### 3.3 Le modèle d'estimation de coût des réservoirs à carburant

Cet exemple porte sur un modèle de coût des réservoirs, dans lequel intervient une variable binaire. Le périmètre technique du produit étudié est le réservoir à carburant équipé de la jauge à carburant et de la pompe d'aspiration. Le réservoir est réalisé en plastique, selon un procédé de soufflage de la matière. Il fait partie d'un sous-ensemble plus vaste du véhicule : le circuit à carburant.

Le réservoir équipé constitue un exemple de fonction du véhicule, dont la conception et la fabrication sont entièrement externalisées et incombent aux fournisseurs équipementiers. Le réservoir est donc entièrement une "pièce œuvrée extérieure". Le prix de revient du réservoir, pour le constructeur, est le prix d'achat. Ce dernier comporte la marge du fournisseur. Elle peut cependant être considérée constante, au regard des conditions de vente correspondant à des volumes constants, que ce soit pour le constructeur automobile comme pour les fournisseurs, dont les produits constituent notre échantillon.

Par ailleurs, il importe de distinguer les réservoirs diesel des réservoirs essence. La discrimination s'explique techniquement par la présence d'une pompe intégrée dans l'ensemble d'aspiration, sur les réservoirs essence, qui n'équipe pas les réservoirs diesel. De plus, l'essence est beaucoup plus volatile que le gazole et les normes de dépollution imposent une perméabilité maximale du réservoir à respecter. Ces contraintes impliquent donc un traitement anti-évaporation spécifique. Notons que le coût du procédé de fluoration dépend essentiellement de la norme de dépollution et marginalement, de la capacité du réservoir.

Ainsi, les réservoirs essence présentent un surcoût par rapport à leurs homologues diesel. De plus, ce surcoût est variable. En effet, des contraintes d'architecture et de conception peuvent impliquer, selon le véhicule, un double puits (un pour le jaugeage et un pour l'aspiration).

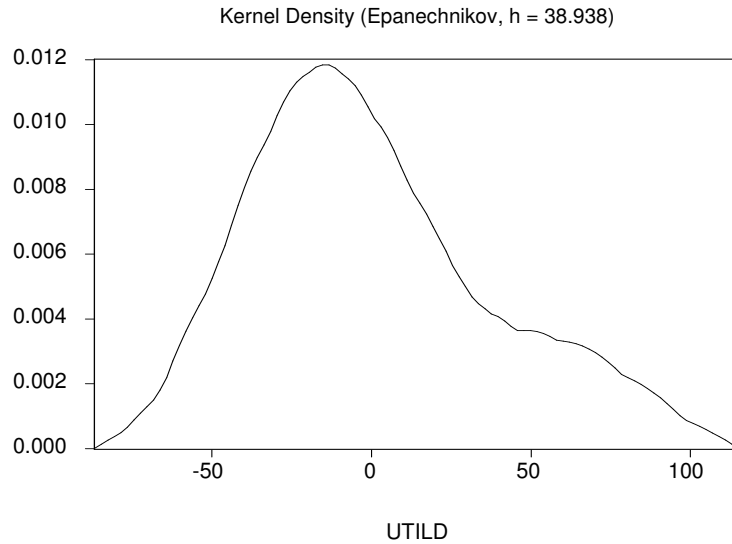
Le modèle économétrique explique le coût en fonction de la capacité du réservoir et d'une variable muette "carburant" qui vaut zéro si le réservoir contient du gazole et un s'il contient de l'essence. Il est estimé par MCO. Le  $R^2=0.95$  laisse apparaître qu'une large part de la variance du coût est expliquée par ces deux variables ; les paramètres estimés des variables capacité et carburant sont significatifs, pour un risque de première espèce de 5 % ; le test de White  $F(3,11)=1.29$  ne permet pas de rejeter l'hypothèse d'homoscédasticité.

Par ailleurs, nous avons effectué un test de Fisher, afin de juger de la validité d'un modèle global, par rapport à deux régressions distinctes, sur chaque type de réservoir (essence et diesel) :  $F(1,11) = 0,54$  qui est inférieur à la valeur critique au seuil 5 %. Ainsi, le test ne permet pas de rejeter l'hypothèse nulle d'une seule régression pour les réservoirs essence et diesel.

La figure 4 présente l'estimateur à noyau de la densité des résidus transformés  $\tilde{u}$ , sans distinction du type de réservoir. La queue de distribution est nettement plus longue vers la

droite, traduisant la présence de résidus positifs importants. En effet, les résidus "extrêmes" correspondant aux observations n° 2 et 8 (réservoirs essence), possèdent des valeurs élevées. D'un point de vue technique, il s'avère que ceux-ci, contrairement aux autres réservoirs essence, sont munis d'un double puits, imposé par les contraintes d'architecture du véhicule.

Figure 4 L'estimateur à noyau de la densité des résidus transformés



– Les procédures de bootstrap des résidus classiques et stratifiés

Nous étudions le processus de construction des échantillons bootstrap qui, en présence d'une variable muette dans la régression, présente certaines spécificités. En effet, la procédure de bootstrap des résidus "classique", telle que nous l'avons présentée dans le paragraphe 1.2, ne respecte pas la structure de l'échantillon initial. Dans le PGD bootstrap, les résidus (essence ou diesel) sont affectés aléatoirement à la capacité d'un réservoir diesel, par exemple. Intuitivement, cette procédure ne semble pas satisfaisante puisque la démarche sous-jacente au bootstrap consiste à générer des échantillons artificiels les plus proches possible de l'échantillon initial. Dès lors, en présence d'une (ou plusieurs) variables explicatives de type qualitatif dans le modèle, nous avons adopté une version adaptée une procédure de bootstrap des résidus "stratifiés".

La construction de l'échantillon bootstrap est détaillée ci-dessous.

Soit  $\tilde{u}$  le vecteur de dimension  $n = 15$ , des résidus transformés de la régression. Ce vecteur est scindé en deux sous-échantillons :  $\tilde{u}_1$  de taille  $n_1 = 8$ , composé des résidus associés aux réservoirs essence et  $\tilde{u}_2$  de taille  $n_2 = 7$ , des 7 résidus associés aux réservoirs diesel. Nous effectuons respectivement 8 (7) tirages aléatoires avec remise dans  $\tilde{u}_1$  ( $\tilde{u}_2$ ) pour constituer  $\tilde{u}_1^*$  ( $\tilde{u}_2^*$ ). La concaténation des deux vecteurs  $\tilde{u}_1^*$  et  $\tilde{u}_2^*$  forme ensuite le vecteur des résidus bootstrap rééchantillonnés "par strates" :  $\tilde{u}^*$  de dimension  $n = 15$ . Les deux procédures de bootstrap des résidus (classique et stratifiée) sont mises en œuvre et comparées lors de la construction d'intervalles de prédiction bootstrap.

– La prédiction de coût

Les prévisions de coût et leurs intervalles, sont calculés pour des réservoirs diesel et essence, de capacité de 60 litres (tableaux 6, 7 et 8).

Tableau 6 Les prévisions MCO et les intervalles de prédiction standard

Réservoir	Prévision MCO	Intervalle de prédiction standard			
		2,5 %	97,5 %	Étendue	Forme
<b>60 l Diesel</b>	260,24	182,74	337,74	157,00	1,00
<b>60 l Essence</b>	507,50	430,37	584,63	154,26	1,00

Unités : francs français

Tableau 7 Les intervalles de prédiction bootstrap des résidus classique

Réservoir	Intervalle de prédiction percentile				Intervalle de prédiction percentile-t			
	2,5 %	97,5 %	Étendue	Forme	2,5 %	97,5 %	Étendue	Forme
<b>60 l Diesel</b>	203,20	342,22	139,02	1,44	203,33	354,27	150,94	1,65
<b>60 l Essence</b>	449,90	588,30	138,40	1,40	451,73	599,26	147,53	1,64

Unités : francs français

Tableau 8 Les intervalles de prédiction bootstrap des résidus stratifiés

Réservoir	Intervalle de prédiction percentile				Intervalle de prédiction percentile-t			
	2,5 %	97,5 %	Étendue	Forme	2,5 %	97,5 %	Étendue	Forme
<b>60 l Diesel</b>	220,52	308,91	88,39	1,22	215,33	312,98	97,65	1,17
<b>60 l Essence</b>	441,50	594,94	153,44	1,32	446,19	619,38	173,19	1,82

Unités : francs français

Les intervalles de prédiction obtenus avec la procédure de bootstrap des résidus classique sont proches, en termes d'étendue, des intervalles standard. Leurs formes sont asymétriques vers les plus fortes valeurs de coûts, ceci de manière identique pour les réservoirs essence et diesel. Ainsi, cette procédure, qui réaffecte de manière aléatoire les résidus, retranscrit l'asymétrie de leur distribution indifféremment sur les deux types de réservoirs. Or, cette asymétrie, causée par les réservoirs essence, n'a pas lieu d'être reportée sur les réservoirs diesel.

La procédure de bootstrap des résidus stratifiés fournit des intervalles de prédiction dont l'étendue est réduite (de l'ordre de 100 francs) pour les réservoirs diesel et plus large (175

francs) pour les réservoirs essence. De plus, ces intervalles sont symétriques pour les réservoirs diesel et fortement asymétrique vers la droite pour les réservoirs essence.

– Le repérage des résidus "extrêmes" dans les répliques bootstrap pour la prédiction

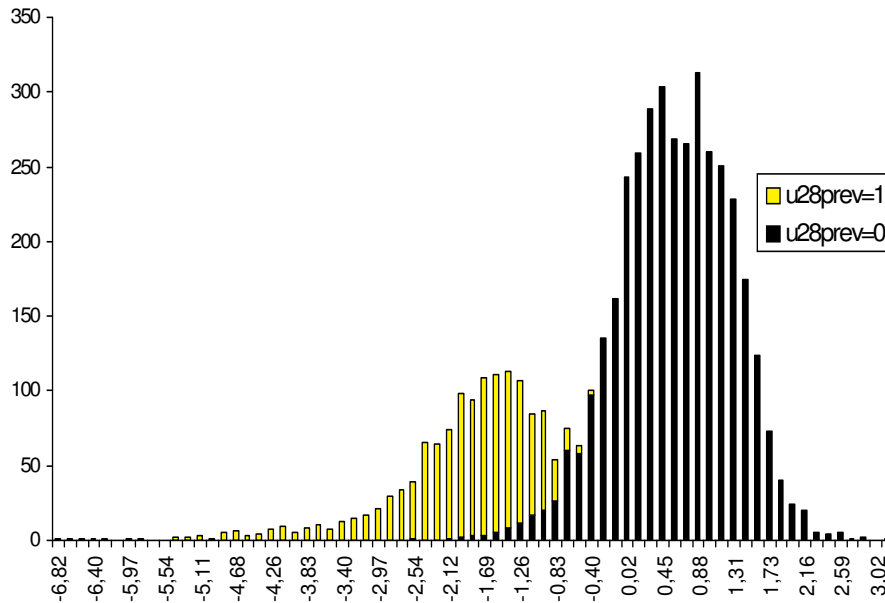
L'impact du tirage des résidus extrêmes sur la distribution de l'erreur de prédiction bootstrap est étudié, en vérifiant pour chaque réplique, si le résidu  $\tilde{u}_f^*$  du modèle théorique bootstrap de la prédiction correspond à  $\hat{u}_2$  ou  $\hat{u}_8$ . Le tableau 9 présente, suivant  $\tilde{u}_f^*$ , les valeurs moyennes et écart-types de la statistique bootstrap de l'erreur de prédiction normée, pour cinq mille répliques bootstrap.

Tableau 9 Les caractéristiques de  $z_f^*$  pour un réservoir de 60 litres

	Réservoir 60 litres	$z_f^*$	$\tilde{u}_f^* \neq (\hat{u}_2$ ou $\hat{u}_8)$	$\tilde{u}_f^* = (\hat{u}_2$ ou $\hat{u}_8)$	Total
<b>Bootstrap des résidus classiques</b>	Diesel	Moyenne	0,256	-1,984	-0,028
		Écart-type	0,793	0,761	1,086
		Effectifs	4365	635	5000
	Essence	Moyenne	0,265	-2,062	-0,033
		Écart-type	0,790	0,895	1,119
		Effectifs	4360	640	5000
<b>Bootstrap des résidus stratifiés</b>	Diesel	Moyenne	0,010	0,002	0,008
		Écart-type	0,733	0,772	0,743
		Effectifs	3723	1277	5000
	Essence	Moyenne	0,603	-1,974	-0,035
		Écart-type	0,685	0,897	1,337
		Effectifs	3763	1237	5000

La figure 5 illustre la distribution empirique bootstrap de l'erreur de lorsque la procédure de bootstrap des résidus est stratifiée. La distribution bootstrap de l'erreur de prédiction normée  $z_f^*$  paraît fortement asymétrique vers les valeurs négatives. Cette asymétrie est due au tirage des résidus extrêmes  $\hat{u}_2$  ou  $\hat{u}_8$ , dans le modèle théorique de prévision bootstrap (cf.. figure 5). En effet, comme  $e_f^* = \hat{y}_f^* - \hat{y}_f - \tilde{u}_f^*$ , si  $\tilde{u}_f^* = \hat{u}_2$  ou  $\hat{u}_8$  (fortes valeurs positives), l'erreur de prédiction bootstrap est fortement négative. Ce phénomène se reproduit pour la distribution de la statistique  $z_f^*$  d'autres prédictions (cf. Juan, 1999)

Figure 5 La distribution de la statistique  $z_f^*$ , pour un réservoir essence de 60 litres (bootstrap des résidus stratifiés)



## CONCLUSION

L'application des méthodes de bootstrap sur les modèles de régression fournit une approximation de la distribution des erreurs de prédiction par leur distribution empirique lorsque celle-ci est inconnue. Le bootstrap est ainsi particulièrement utile lorsque les échantillons de données sont de petite taille et qu'il n'est pas possible de formuler l'hypothèse d'une distribution gaussienne du terme d'erreur. Le nombre de répliques peut être déterminé à partir des coefficients de variation de l'étendue de l'intervalle de confiance des coefficients ou de l'intervalle de prédiction lorsque ceux-ci deviennent peu variant.

La mise en œuvre du bootstrap invite à privilégier le calcul des paramètres estimés en utilisant la pseudo-inverse de la matrice des variables explicatives. Cette méthode de calcul s'avère robuste en présence de multicollinéarité. Elle est également performante en terme de temps de calcul car il est suffisant de réduire les données qui sont alors normées à un avant de calculer l'estimateur et non pas de les centrer et les réduire comme pour une inversion de matrice.

Un algorithme modifié permet de trier rapidement la distribution empirique de la statistique d'intérêt. Seules les queues de distribution sont triées et les éléments sont comparés aux valeurs extrêmes de celles-ci qui correspondent aux fractiles retenus.

La prévision des coûts au stade des avant-projets dans l'industrie automobile soulève des difficultés liées à la faible taille des échantillons de données et à l'asymétrie qui caractérise des distributions de coût. Deux applications permettent d'apprécier l'apport du bootstrap.

Le premier exemple développé illustre l'utilisation des techniques de bootstrap sur un modèle simplifié de coût d'un moteur. L'analyse des intervalles de prédiction montre que le bootstrap permet de retranscrire l'asymétrie de la distribution des résidus dans les intervalles de prédiction. En effet, ces derniers sont décalés, par rapport aux intervalles standard, vers les plus fortes valeurs de coûts et autorisent ainsi, pour la prédiction du coût d'un nouveau moteur, des valeurs plus élevées. L'utilisation des techniques de bootstrap permet donc une meilleure retranscription de l'information contenue dans l'échantillon initial, pour les intervalles de prédiction.

Le second exemple a permis d'exposer une utilisation des techniques de bootstrap adaptée, dans le cas d'une modélisation en présence de variable muette. Cette méthode permet de construire des intervalles de prédiction symétriques pour les réservoirs diesel et asymétriques vers les plus fortes valeurs de coûts pour les réservoirs essence. Ceci résulte de l'asymétrie, à la fois de la distribution de l'erreur de prédiction et de celle du surcoût lié au carburant essence. Ainsi, l'utilisation de la procédure de bootstrap des résidus stratifiés permet, lorsque l'information n'est pas disponible, de prendre en compte dans les intervalles de confiance et de prédiction, des surcoûts éventuels pour les réservoirs essence, imposés par les contraintes d'architecture du véhicule.

Des développements dans l'utilisation du bootstrap sont envisagés, notamment la prise en compte d'éventuelles non-linéarités dans la spécification des fonctions de coût. En effet, bien que ceci n'ait pas été mis en évidence dans les exemples traités, il faut envisager la non-linéarité de la relation comme une alternative à un modèle linéaire avec terme d'erreur asymétrique. L'utilisation du bootstrap concerne alors les tests de spécifications ainsi que l'estimation des paramètres du modèle.

## Références

- Belsley D., Kuh E., Welsch R., 1980, “ *Regression diagnostics, identifying influential data and sources of collinearity* ”, ed. John Wiley and Sons, New-York
- Belsley D., 1984, “ Demeaning conditioning diagnostic through centering ”, *The American Statistician*, Vol. 38, n°2, p 73-93
- Breiman L. , 1992, “The little bootstrap and other methods for dimensionality selection in regression : X-fixed prediction error ”, *Journal of the American Statistical Association*, Vol. 87, p. 738-754.
- Booth J.G., Sarkar S., 1998, “Monte-Carlo approximation of bootstrap variances”, *The American Statistician*, Vol. 52, n°4, p 354-357
- Davidson R., McKinnon J., 1998, “Bootstrap tests : how many bootstraps ? ”, Document de travail, Université de la Méditerranée, GREQAM
- Efron B. , 1979, “ Bootstrap methods : another look at the jackknife ”, *Annals of Statistics*, Vol. 7, p. 1-26.
- Efron B., Tibshirani R.J. , 1993, “ *An introduction to the bootstrap* ”, ed. Chapman and Hall, New-York
- Erkel-Rousse H., 1995 ,“ Détection de la multicollinéarité dans un modèle linéaire ordinaire : quelques éléments pour un usage averti des indicateurs de Belsley, Kuh et Welsch ”, *Revue de statistique appliquée*, Vol. 18, n°4, p 19-42
- Flachaire E., 1998, “ les méthodes du bootstrap et l’inférence robuste à l’hétéroscédasticité“, Thèse de doctorat, Université de la Méditerranée, GREQAM
- Forsythe G.E., Malcolm M.A., Moler C.B., 1977, “ *Computer methods for mathematical computations* ”, Prentice-Hall
- Freedman D. A. , 1981, “ Bootstrapping regression models ”, *The Annals of Statistics*, Vol. 9, p. 1218-1228.
- Hall P., 1992, “ *The bootstrap and edgeworth expansion* ”, ed. Springer Verlag, New-York
- Horowitz J.L., 1997, “ Bootstrap methods in econometrics : theory and numerical performance”, “ *Advances in economics and econometrics : theory and application* ”, Vol. 3, Cambridge University Press, p.188-222.
- Jeong J., Maddala G.S., 1993, “ A perspective on application of bootstrap methods in econometrics ”, *Handbook of Statistics*, Vol. 11, Amsterdam : North-Holland, p. 573-610.
- Juan S., 1999, “ Les modélisation économétriques d’estimation de coût dans l’industrie automobile: l’apport des techniques de bootstrap“, Thèse de doctorat, ENSPM-Université de Bourgogne
- Lantz F., 1983, “ *Mise en œuvre de la régression linéaire : calcul et stabilité de la matrice (X’X) inverse*”, Les cahiers du 3<sup>e</sup> cycle économétrie , Université de Paris X, p. 34-47
- McCullough B.D., 1996, “ Estimating forecast intervals when the exogenous variable is stochastic ”, *Journal of Forecasting*, Vol. 15, p. 293-304.
- McCullough B.D., Vinod H., 1993, “ Implementing the single bootstrap : some computational considerations”, *Computational Economics*, Vol. 6, p. 1-15.
- Seaks T., 1972, “ Computer algorithms - syminv : an algorithm for the inversion of a positive definite matrix by the Cholesky decomposition”, *Econometrica*, Vol. 40, n°5
- Simar L., 1992, “ Estimating efficiencies from frontier models with panel data : a comparison of parametric, non parametric and semi-parametric methods with bootstrapping ”, *The Journal of Productivity Analysis*, Vol.3, p. 171-203.

- Stine R.A ., 1985, “ Bootstrap prediction intervals for regression ”, *Journal of The American Statistical Association*, Vol. 80, p. 1026-1031.
- Veall M.R., 1998, “Applications of the bootstrap”, *Handbook of Applied Economic Statistics*, Vol. 155, ed. Aman U., Giles D.E.A.
- Verboven F., 1996, “ International price discrimination in the european car market”, *Rand Journal of Economics*, Vol. 27, p 240-268
- Vinod H., 1993, “Bootstrap methods : applications in econometrics”, *Handbook of Statistics*, Vol. 11, North-Holland, p.629-661.



## Annexe 1 - Principe du bootstrap et application au modèle de régression

### A.1 - Le principe

Le principe du bootstrap consiste, en répétant un grand nombre de fois le rééchantillonnage dans les données d'origine, à construire la fonction de répartition empirique bootstrap d'une statistique d'intérêt. Cette dernière approche alors de manière satisfaisante la vraie distribution de la statistique qui, elle, est inconnue. Le processus de construction de la fonction de répartition empirique bootstrap d'un estimateur est détaillé ci-dessous.

Soit un échantillon i.i.d.  $\{y_i\}_{i=1}^n$ , d'une variable aléatoire  $y$  de loi  $F$  inconnue, nous cherchons à déterminer la loi  $\hat{\theta}(F)$  d'un estimateur  $\hat{\theta}$  d'un paramètre  $\theta$  de  $F$ . L'objectif poursuivi est donc de construire une table statistique<sup>10</sup> de valeurs approchées de  $\hat{\theta}(F)$ . Pour ce faire, nous disposons des  $n$  observations de l'échantillon, donc de la fonction de répartition empirique  $\hat{F}_n$ .

En théorie, il est possible de construire une table de  $\hat{\theta}(\hat{F}_n) : T$ , conditionnelle aux  $\{y_i\}_{i=1}^n$ , en calculant  $\hat{\theta}$  sur chacun des  $n$ -échantillons tirés avec remise dans  $\{y_i\}_{i=1}^n$ . Cependant, il existe  $n^n$  tels échantillons et cette procédure ne peut être mise en œuvre que lorsque  $n$  est très petit. En pratique, nous allons tirer  $B$  échantillons ( $b = 1, \dots, B$ ), nommés échantillons bootstrap, pour construire une table extraite de  $T$ . Sur chacun des échantillons bootstrap, nous calculons la valeur  $\hat{\theta}^*(b)$  de la statistique  $\hat{\theta}$ . La table bootstrap de  $\hat{\theta}^*$  est donc une sous-table de  $T$ , conditionnelle à l'échantillon de données. Cette notion de dépendance aux données est importante, pour la compréhension du processus bootstrap. La table bootstrap ne s'applique en effet, que pour l'échantillon initial. Pour un nouvel échantillon, il est donc nécessaire de construire une nouvelle table bootstrap, propre à cet échantillon.

D'après le théorème de Glivenko-Cantelli (G-C), pour un échantillon de variables aléatoires i.i.d. de loi  $F$  inconnue, lorsque la taille de l'échantillon  $n$  tend vers l'infini, la fonction de répartition empirique de l'échantillon converge uniformément presque sûrement vers la loi  $F$ . Ainsi, la table bootstrap de  $\hat{\theta}^*$ , conditionnelle aux  $\{y_i\}_{i=1}^n$ , fournit une bonne approximation de la loi  $\hat{\theta}(F)$ .

### A.2 - Les méthodes bootstrap sur les modèles de regression

Le modèle de régression linéaire multiple est noté :

$$Y = X\beta + u \quad (\text{A-1})$$

---

<sup>10</sup> Par « table statistique », nous désignons une version empirique (sur les  $B$  répliques bootstrap) de la distribution d'échantillonnage de  $\hat{\theta}$ .

où  $Y$  est un vecteur  $(n,1)$ ,  $X$  une matrice  $(n,p)$ ,  $\beta$  le vecteur des coefficients à estimer  $(p,1)$  et  $u$  le vecteur des erreurs aléatoires  $(n,1)$ . Un rang d'observations  $i$  ( $i = 1, \dots, n$ ) de la matrice  $X$ , correspondant à une ligne, est noté  $X_i$   $(1,p)$ . Les paramètres estimés par la méthode des moindres carrés ordinaires (MCO)  $\hat{\beta}$  et les résidus  $\hat{u}$  sont définis comme:  $\hat{\beta} = (X^T X)^{-1} X^T Y$  et  $\hat{u} = Y - X\hat{\beta}$ .

### – Le bootstrap des résidus

Le modèle théorique bootstrap est le suivant :

$$Y^* = X\hat{\beta} + u^* \quad (\text{A-2})$$

où  $\hat{\beta}$  est l'estimateur MCO et  $u^*$  est un terme aléatoire issu des résidus  $\hat{u}$  de la régression initiale, dont nous décrivons la construction ci-dessous.

L'application de la procédure bootstrap consiste à répéter  $B$  fois les étapes suivantes :

1) A chaque itération  $b$  ( $b=1, \dots, B$ ), un échantillon  $\{y_i^*\}_{i=1}^n$ , de dimension  $(n,1)$ , est constitué à partir du modèle bootstrap (A-2). Nous disposons alors d'un nouveau couple  $(Y^*, X)$  à partir duquel nous pouvons réaliser une estimation des paramètres de la régression.

2)

Les résidus MCO étant plus petits que les erreurs qu'ils estiment, une transformation est nécessaire pour élaborer le terme aléatoire du modèle théorique bootstrap. Ainsi, à la suite de Freedman (1981), ce dernier est construit avec les résidus transformés suivants ( $\hat{u}_i$  est divisé par un facteur proportionnel à la racine de sa variance) qui sont de même norme que que les termes erreurs  $u_i$  :

$$\tilde{u}_i = \frac{\hat{u}_i}{\sqrt{(1-h_i)}} - \frac{1}{n} \sum_{s=1}^n \frac{\hat{u}_s}{\sqrt{(1-h_s)}}$$

Le modèle théorique bootstrap est donc le suivant :

$$y_i^*(b) = X_i \hat{\beta} + \tilde{u}_i^*(b), \quad i = 1, \dots, n \quad (\text{A-3})$$

Où  $\tilde{u}_i^*(b)$  est rééchantillonné à partir des  $\tilde{u}_i$ . Les nouvelles variables dépendantes bootstrappées  $y_i^*(b)$  sont donc construites à partir des valeurs calculées  $\hat{y}_i$  et des  $\tilde{u}_i^*(b)$  :  $y_i^*(b) = \hat{y}_i + \tilde{u}_i^*(b)$ .

3) La procédure d'estimation par MCO est appliquée sur le modèle de régression (A-3) afin d'obtenir l'estimateur bootstrap. Pour le  $b$ -ième échantillon, ce dernier s'écrit :

$$\hat{\beta}^*(b) = (X^T X)^{-1} X^T Y^*(b) \quad (\text{A-4})$$

Les étapes 1) et 2) sont répétées  $B$  fois ( $b=1, \dots, B$ ). Les fonctions de répartition empirique bootstrap de  $\hat{\beta}^*$  et des statistiques issues de  $\hat{\beta}^*$  sont ensuite construites.

Lors de la mise en œuvre d'une procédure de bootstrap sur les résidus, les variables explicatives  $X$  sont considérées comme fixées. Cette procédure est donc valide si les  $X$  sont réellement fixés et si les erreurs vérifient les hypothèses classiques des MCO<sup>11</sup>. Par conséquent, elle n'est pas correcte si ces dernières sont hétéroscédastiques. En appliquant la procédure de bootstrap des résidus pour construire le modèle théorique bootstrap, différents termes erreurs sont associés à différentes variables explicatives. Ainsi, la procédure bootstrap de rééchantillonnage des résidus n'est pas en mesure de respecter cette relation. Dans de tels cas, le modèle théorique bootstrap est construit suivant une procédure différente, nommée "bootstrap par paires".

### – Le bootstrap par paires

Cette seconde approche bootstrap des modèles de régression consiste à rééchantillonner directement dans les données d'origine, à partir des paires  $(y_i, X_i)$ . Notons cependant que le tirage simultané de  $(y_i, X_i)$  introduit une corrélation entre les régresseurs et les erreurs du processus générateur de données (PGD) bootstrap. Ainsi, le bootstrap par paires, sous cette forme<sup>12</sup> simple, ne respecte pas l'hypothèse d'exogénéité des régresseurs dans le PGD bootstrap.

L'application de la procédure bootstrap par paires consiste à répéter  $B$  fois les étapes suivantes :

1) A chaque itération  $b$  ( $b=1, \dots, B$ ), le vecteur  $Y^*$  et la matrice des variables explicatives  $X^*$  sont construits, en effectuant  $n$  tirages aléatoires avec remise<sup>13</sup> de paires  $(y_i, X_i)$ , dans l'échantillon d'origine. Ainsi, si le terme erreur  $u_i$  associé à  $X_i$  a une grande variance, la relation sera préservée dans l'échantillon bootstrap.

2) Une estimation par MCO des coefficients du modèle de régression bootstrap est ensuite réalisée :

$$\hat{\beta}^*(b) = (X^{*T}(b)X^*(b))^{-1} X^{*T}(b)Y^*(b) \quad (\text{A-5})$$

Notons, à la différence de la procédure bootstrap des résidus, que la matrice des variables explicatives  $X^*(b)$  est différente, à chaque itération  $b$ .

Les  $B$  répliques  $\hat{\beta}^*$  fournissent alors la fonction de répartition empirique bootstrap. Ainsi, les  $B$  répliques bootstrap indépendantes, obtenues suivant les procédures de bootstrap présentées ci-dessus, fournissent un échantillon aléatoire des  $\hat{\beta}^*$  qui est utilisé

<sup>11</sup> Les erreurs ont une espérance mathématique nulle, sont homoscédastiques et non-autocorrélées.

<sup>12</sup> Le *wild bootstrap* est une méthode adaptée en présence d'hétéroscédasticité, qui respecte l'hypothèse d'exogénéité des régresseurs.

<sup>13</sup> Notons que Freedman (1981) envisage des échantillons bootstrap de taille  $m$  différente de  $n$ .

pour estimer la distribution bootstrap de  $\hat{\beta}$ . Cette dernière permet alors la construction des intervalles de confiance bootstrap des paramètres du modèle de régression.

### A.3 - Construction des intervalles de confiance bootstrap

Nous rappelons brièvement la forme des intervalles de confiance standard des coefficients de la régression, qui dépendent de l'hypothèse de Normalité des termes erreurs, avant de présenter les intervalles de confiance bootstrap, qui dépendent uniquement des données. Les développements sont effectués pour l'élément  $j$  du vecteur des paramètres  $\beta$ .

Soit la variable aléatoire  $z_j$  définie comme  $z_j = \frac{\hat{\beta}_j - \beta_j}{s(\hat{\beta}_j)}$ , l'intervalle de confiance standard

de  $\beta_j$  découle de l'hypothèse selon laquelle  $z_j$  est distribuée selon une loi de Student à  $n-p$  degrés de liberté. Ainsi, pour un niveau de confiance  $(1-2\alpha)$ , l'intervalle de confiance standard prend alors la forme suivante :

$$\left[ \hat{\beta}_j - s(\hat{\beta}_j) \cdot t_{(1-\alpha), n-p}, \hat{\beta}_j + s(\hat{\beta}_j) \cdot t_{(\alpha), n-p} \right] \quad (\text{A-6})$$

où  $t$  sont les valeurs des quantiles  $(\alpha)$  et  $(1-\alpha)$  de la distribution de Student à  $n-p$  degrés de liberté.

Les intervalles de confiance bootstrap sont maintenant présentés, au travers de leurs deux principales méthodes de construction : l'approche percentile et percentile-t

#### – La méthode percentile

Cette méthode, basée uniquement sur les estimations bootstrap, est la méthode la plus simple d'obtention d'intervalles de confiance. La première étape consiste à trier les estimations  $\hat{\beta}_j^*(b)$  par ordre croissant, pour les  $b=1, \dots, B$  réplifications. Notons  $\hat{\beta}_j^*(1)$  le plus petit et  $\hat{\beta}_j^*(B)$  le plus grand.

Pour un niveau  $(1-2\alpha)$ , l'intervalle de confiance percentile pour le paramètre  $\beta_j$  est alors donné par :

$$\left[ \hat{\beta}_j^*(\alpha B), \hat{\beta}_j^*((1-\alpha)B) \right] \quad (\text{A-7})$$

$\hat{\beta}_j^*(\alpha B)$  représente la  $\alpha B$ -ième valeur (respectivement  $\hat{\beta}_j^*((1-\alpha)B)$  la  $(1-\alpha)B$ -ième valeur) de la liste ordonnée des  $B$  réplifications bootstrap. Les valeurs seuils sont donc choisies telles que  $\alpha$  % des réplifications ont fourni des  $\hat{\beta}_j^*$  plus petits (grands) que la borne inférieure (supérieure) de l'intervalle de confiance percentile.

#### – La méthode percentile-t

Cette seconde approche de construction d'intervalle de confiance suit une démarche très proche de celle utilisée pour élaborer l'intervalle de confiance standard de  $\beta_j$ . La procédure bootstrap percentile-t consiste à estimer la fonction de répartition de  $z_j$  directement à partir

des données. Cela revient à construire une table statistique à partir de la fonction de répartition empirique des  $B$  répliques bootstrap  $z_j^*$ . Cette table est nommée table bootstrap. Les  $z_j^*$  sont définies comme :

$$z_j^* = \frac{\hat{\beta}_j^* - \hat{\beta}_j}{s^*(\hat{\beta}_j^*)} \quad (\text{A-8})$$

Notons, en comparaison avec la méthode percentile, un calcul supplémentaire dans cette approche. En effet, pour chacune des répliques bootstrap, il est nécessaire de calculer l'écart-type estimé bootstrap  $s^*(\hat{\beta}_j^*)$ .

Soit  $\hat{F}_{z_j^*}$  la fonction de répartition empirique des  $z_j^*$ , le fractile à  $\alpha$  %,  $\hat{F}_{z_j^*}^{-1}(\alpha)$ , est estimé par la valeur  $\hat{t}^{(\alpha)}$  telle que  $\# \{z_j^*(b) \leq \hat{t}^{(\alpha)}\} / B = \alpha$ .

Finalement, l'intervalle de confiance percentile-t pour  $\beta_j$  s'écrit :

$$\left[ \hat{\beta}_j - s(\hat{\beta}_j) \cdot \hat{t}^{(1-\alpha)}, \hat{\beta}_j - s(\hat{\beta}_j) \cdot \hat{t}^{(\alpha)} \right] \quad (\text{A-9})$$

Ainsi, l'intervalle de confiance percentile-t est l'analogie bootstrap de l'intervalle de confiance standard.

En résumé, l'intervalle de confiance percentile-t substitue, aux valeurs critiques de la loi de Student utilisées dans l'intervalle standard, les valeurs seuils de la table bootstrap. Notons que ces dernières peuvent être très différentes. Cette différence est d'autant plus importante que la distribution (inconnue) des erreurs est éloignée de la loi Normale. De plus, nous remarquons que les valeurs des quantiles ( $\alpha$ ) et  $(1-\alpha)$  de la distribution de Student, symétriques par nature, entraînent directement la symétrie de l'intervalle de confiance standard autour de l'estimation  $\hat{\beta}_j$ . Par opposition, les valeurs  $\hat{t}^{(\alpha)}$  et  $\hat{t}^{(1-\alpha)}$  de la table bootstrap peuvent être asymétriques et permettent alors des intervalles de confiance asymétriques autour de  $\hat{\beta}_j$ . Cette prise en compte d'une possible asymétrie constitue un avantage important des intervalles de confiance bootstrap.

## **Déjà parus**

### **CEG-1. D. PERRUCHET, J.-P. CUEILLE,**

Compagnies pétrolières internationales : intégration verticale et niveau de risque.  
Novembre 1990

### **CEG-2. C. BARRET, P. CHOLLET,**

Canadian gas exports: modeling a market in disequilibrium.  
Juin 1990

### **CEG-3. J.-P. FAVENNEC, V. PREVOT,**

Raffinage et environnement.  
Janvier 1991

### **CEG-4. D. BABUSIAUX,**

Note sur le choix des investissements en présence de rationnement du capital.  
Janvier 1990

### **CEG-5. J.-L. KARNIK,**

Les résultats financiers des sociétés de raffinage distribution en France 1978-89.  
Mars 1991

### **CEG-6. I. CADORET, P. RENOU,**

Élasticités et substitutions énergétiques : difficultés méthodologiques.  
Avril 1991

### **CEG-7. I. CADORET, J.-L. KARNIK,**

Modélisation de la demande de gaz naturel dans le secteur domestique : France, Italie, Royaume-Uni 1978-1989.  
Juillet 1991

### **CEG-8. J.-M. BREUIL,**

Émissions de SO<sub>2</sub> dans l'industrie française : une approche technico-économique.  
Septembre 1991

### **CEG-9. A. FAUVEAU, P. CHOLLET, F. LANTZ,**

Changements structurels dans un modèle économétrique de demande de carburant.  
Octobre 1991

### **CEG-10. P. RENOU,**

Modélisation des substitutions énergétiques dans les pays de l'OCDE.  
Décembre 1991

### **CEG-11. E. DELAFOSSE,**

Marchés gaziers du Sud-Est asiatique : évolutions et enseignements.  
Juin 1992

### **CEG-12. F. LANTZ, C. IOANNIDIS,**

Analysis of the French gasoline market since the deregulation of prices.  
Juillet 1992

### **CEG-13. K. FAID,**

Analysis of the American oil futures market.  
Décembre 1992

### **CEG-14. S. NACHET,**

La réglementation internationale pour la prévention et l'indemnisation des pollutions maritimes par les hydrocarbures.  
Mars 1993

**CEG-15. J.-L. KARNIK, R. BAKER, D. PERRUCHET,**

Les compagnies pétrolières : 1973-1993, vingt ans après.

Juillet 1993

**CEG-16. N. ALBA-SAUNAL,**

Environnement et élasticités de substitution dans l'industrie ; méthodes et interrogations pour l'avenir.

Septembre 1993

**CEG-17. E. DELAFOSSE,**

Pays en développement et enjeux gaziers : prendre en compte les contraintes d'accès aux ressources locales.

Octobre 1993

**CEG-18. J.P. FAVENNEC, D. BABUSIAUX\*,**

L'industrie du raffinage dans le Golfe arabe, en Asie et en Europe : comparaison et interdépendance.

Octobre 1993

**CEG-19. S. FURLAN,**

L'apport de la théorie économique à la définition d'externalité.

Juin 1994

**CEG-20. M. CADREN,**

Analyse économétrique de l'intégration européenne des produits pétroliers : le marché du diesel en Allemagne et en France.

Novembre 1994

**CEG-21. J.L. KARNIK, J. MASSERON\*,**

L'impact du progrès technique sur l'industrie du pétrole.

Janvier 1995

**CEG-22. J.P. FAVENNEC, D. BABUSIAUX,**

L'avenir de l'industrie du raffinage.

Janvier 1995

**CEG- 23. D. BABUSIAUX, S. YAFIL\*,**

Relations entre taux de rentabilité interne et taux de rendement comptable.

Mai 1995

**CEG-24. D. BABUSIAUX, J. JAYLET\*,**

Calculs de rentabilité et mode de financement des investissements, vers une nouvelle méthode ?

Juin 1996

**CEG-25. J.P. CUEILLE, J. MASSERON\*,**

Coûts de production des énergies fossiles : situation actuelle et perspectives.

Juillet 1996

**CEG-26. J.P. CUEILLE, E. JOURDAIN,**

Réductions des externalités : impacts du progrès technique et de l'amélioration de l'efficacité énergétique.

Janvier 1997

**CEG-27. J.P. CUEILLE, E. DOS SANTOS,**

Approche évolutionniste de la compétitivité des activités amont de la filière pétrolière dans une perspective de long terme.

Février 1997

**CEG-28. C. BAUDOUIN, J.P. FAVENNEC,**

Marges et perspectives du raffinage.

Avril 1997

**CEG-29. P. COUSSY, S. FURLAN, E. JOURDAIN, G. LANDRIEU, J.V. SPADARO, A. RABL,**  
Tentative d'évaluation monétaire des coûts externes liés à la pollution automobile : difficultés  
méthodologiques et étude de cas.  
Février 1998

**CEG-30. J.P. INDJEHAGOPIAN, F. LANTZ, V. SIMON,**  
Dynamique des prix sur le marché des fiouls domestiques en Europe.  
Octobre 1998

**CEG-31. A. PIERRU, A. MAURO,**  
Actions et obligations : des options qui s'ignorent.  
Janvier 1999

**CEG-32. V. LEPEZ, G. MANDONNET,**  
Problèmes de robustesse dans l'estimation des réserves ultimes de pétrole conventionnel.  
Mars 1999

**CEG-33. J. P. FAVENNEC, P. COPINSCHI,**  
L'amont pétrolier en Afrique de l'Ouest, état des lieux  
Octobre 1999

**CEG-34. D. BABUSIAUX,**  
Mondialisation et formes de concurrence sur les grands marchés de matières premières énergétiques : le  
pétrole.  
Novembre 1999

**CEG-35. D. RILEY,**  
The Euro  
Février 2000

**CEG-36. et 36bis D. BABUSIAUX, A. PIERRU\*,**  
Calculs de rentabilité et mode de financement des projets d'investissements : propositions méthodologiques.  
Avril 2000 et septembre 2000

**CEG-37. P. ALBA, O. RECH,**  
Peut-on améliorer les prévisions énergétiques ?  
Mai 2000

**CEG-38. J.P. FAVENNEC, D. BABUSIAUX,**  
Quel futur pour le prix du brut ?  
Septembre 2000

\* une version anglaise de cet article est disponible sur demande