

# La régression quantile en pratique

Xavier D'Haultfœuille\* et Pauline Givord\*\*

L'usage des régressions quantiles s'est beaucoup répandu au cours de la dernière décennie. Celles-ci reposent sur un principe proche de celui de la régression linéaire classique. De même que cette dernière se fonde sur une modélisation linéaire de l'espérance conditionnelle de la variable d'intérêt en fonction de ses déterminants, les régressions quantiles consistent à supposer que les quantiles conditionnels de cette variable d'intérêt sont linéaires. Elles fournissent cependant une description plus riche que les régressions linéaires, puisqu'on peut ainsi étudier l'ensemble de la distribution conditionnelle de la variable d'intérêt et non seulement la moyenne de celle-ci. Cette analyse est particulièrement intéressante pour les mesures d'évaluation des politiques publiques : un programme peut avoir un effet moyen limité, mais permettre d'augmenter suffisamment les niveaux les plus faibles de la variable d'intérêt pour que son implémentation soit souhaitable. Les régressions quantiles permettent également de décrire les déterminants des évolutions des inégalités de revenu. En outre, elles sont parfois plus adaptées pour certains types de données (variables censurées ou tronquées, présence de valeurs extrêmes, modèles non linéaires...).

Les régressions quantiles peuvent être aujourd'hui effectuées aisément avec de nombreux logiciels statistiques. Cet article rappelle les principes statistiques sous-jacents à cette modélisation, ainsi que des extensions qui ont été développées pour répondre au problème, classique en économétrie, de l'endogénéité de certaines variables explicatives (données de panel, variables instrumentales...). Il fournit également un guide d'interprétation des résultats d'une régression quantile, dont l'analyse est peut-être moins intuitive que celle d'une régression linéaire. Pour bien comprendre l'utilisation qui peut en être faite, deux applications concrètes sont présentées à titre d'illustration.

## Rappel :

Les jugements et opinions exprimés par les auteurs n'engagent qu'eux mêmes, et non les institutions auxquelles ils appartiennent, ni a fortiori l'Insee.

\* Crest. xavier.dhaultfoeuille@ensae.fr

\*\* Insee, Département des Méthodologies Statistiques, Direction de la Méthodologie et de la Coordination Statistique Internationale. pauline.givord@insee.fr

Les auteurs remercient les très nombreuses personnes qui ont contribué par leurs commentaires sur les premières versions successives de cet article à l'améliorer significativement, et plus particulièrement Cédric Afsa, Pascale Breuil, Elise Coudin, Jean-Michel Floch, Marine Guillem, Jérôme Le, Simon Quantin et Olivier Sautory, ainsi que deux relecteurs anonymes de la revue. Ils restent seuls responsables des erreurs et approximations qui pourraient demeurer dans cet article.

Beaucoup d'études économiques empiriques se concentrent sur l'observation ou la modélisation de la moyenne et cette focalisation est souvent critiquée. La moyenne apporte une information essentielle mais néanmoins limitée. Le revenu moyen n'informe pas sur la répartition plus ou moins inégale de ces revenus dans la population. L'une des préconisations du rapport Stiglitz-Sen-Fitoussi appelait ainsi à sortir de la dictature de la moyenne. Cette préconisation visait surtout à l'enrichissement des constats descriptifs sur le revenu et les niveaux de vie, mais elle peut s'étendre aux travaux de type plus analytiques qui s'intéressent à leurs déterminants. Raisonner en termes de distributions est également utile dans d'autres contextes que la mesure des inégalités. Par exemple, pour évaluer l'effet d'une politique publique, il est souvent pertinent d'aller au-delà de ses effets moyens. Il peut ainsi être souhaitable de mettre en œuvre une politique éducative qui permet de réduire la proportion d'élèves en grande difficulté, même si elle n'a qu'un effet négligeable sur le niveau moyen de l'ensemble des élèves.

Une deuxième limite du recours à la moyenne est d'ordre plus technique : c'est le fait qu'elle s'avère parfois difficile à modéliser. Cela peut être le cas en présence de valeurs extrêmes ou aberrantes (dues par exemple à des erreurs de mesures), auxquelles la moyenne est bien plus sensible que les quantiles. Lorsque la distribution de la variable d'intérêt est très étalée, ce qui est par exemple le cas des revenus, la moyenne pourra beaucoup varier en fonction de l'échantillon utilisé. L'estimation de la moyenne est également compromise en présence de données censurées, c'est-à-dire lorsqu'on n'observe la variable d'intérêt qu'au-delà ou en deçà d'un seuil fixe. Ainsi, pour des raisons de confidentialité, les données individuelles de revenus sont parfois diffusées en écrétant ceux qui sont supérieurs à un certain niveau. Il n'est pas possible d'estimer la moyenne d'une variable censurée de la sorte, sauf à faire des hypothèses paramétriques sur la distribution de cette variable au-dessus du seuil. En revanche, en-dessous de ce niveau, les quantiles de la variable censurée coïncident avec ceux de la variable d'intérêt (cf. Buchinsky, 1994, pour une application à l'évolution des revenus aux États-Unis).

La régression quantile est l'un des outils dont dispose le statisticien pour répondre à ces limites inhérentes à la moyenne. Elle permet d'avoir une description plus précise de la distribution d'une variable d'intérêt conditionnelle à ses déterminants, comparativement à la régression linéaire

qui se focalise sur la moyenne conditionnelle. Si son principe est ancien<sup>1</sup>, elle a connu récemment un regain d'intérêt. Un ensemble de procédures préprogrammées en font aujourd'hui un outil simple d'utilisation. Pour ne donner que quelques exemples de travaux récents sur données françaises, la régression quantile a été mobilisée pour étudier l'évolution des inégalités des salaires en population générale (Charnoz *et al.*, 2011), pour des comparaisons des distributions des salaires masculins et féminins selon le secteur d'activité (Etienne et Narcy, 2010), ou encore des analyses plus ciblées sur l'éventail de rémunération des médecins (Samson, 2006 ; Dumontet et Franc, 2014). Hors du thème des salaires, Biscourp *et al.* (2013) l'utilisent pour analyser comment la dispersion des prix des biens de grande consommation a réagi aux réformes des relations commerciales qui ont marqué les années 1990 et 2000. La méthode peut même parfois servir à des travaux de type plus macroéconomiques : Cornec (2014) l'utilise pour modéliser les intervalles de confiance entourant les prévisions économiques de court-terme.

Dans ce contexte, l'objet du présent article est de préciser les principes de cette méthode et l'utilisation qui peut en être faite, plus particulièrement pour l'analyse microéconomique. Il peut se lire de plusieurs manières. Les lecteurs souhaitant simplement disposer d'un guide d'interprétation pour les travaux qui mettent en œuvre ce type de méthode tireront surtout profit de la première section qui présente les principaux intérêts de la régression quantile, ainsi que de la section finale qui détaille deux applications pratiques. Un lecteur néophyte souhaitant mettre en œuvre la régression quantile sous sa forme classique complètera cette lecture par la deuxième section qui rappelle les principes d'estimation sous-jacents et les propriétés statistiques des estimateurs. Enfin, un lecteur plus averti trouvera dans la troisième section plusieurs extensions, en particulier sur la prise en compte de l'endogénéité et les modèles non-linéaires. Les parties les plus techniques et pouvant être sautées en première lecture sont identifiées par un astérisque.

On ajoutera que le présent article n'aborde pas les questions de mise en œuvre informatique mais plusieurs exemples sont fournis dans Givord et D'Haultfœuille (2013), dont cet article est issu. Enfin, le lecteur intéressé pourra trouver

1. La méthode des moindres déviations absolues (Least Absolute Deviation, LAD), qui revient à s'intéresser à la médiane plutôt qu'à la moyenne comme dans les moindres carrés ordinaires, est aussi ancienne que la méthode des moindres carrés ordinaires et remonte à Boscovitch, Laplace et Gauss.

d'autres présentations de la méthode, en anglais, dans Koenker et Hallock (2001), Cade et Noon (2003) ou encore Koenker (2005).

## Les apports de la régression quantile

Pour poser les notations, nous nous intéressons à une variable aléatoire  $Y$ , de fonction de répartition  $\tau$  définie par  $F_Y(y) = P(Y \leq y)$ . Rappelons que le quantile d'ordre  $\tau$  est généralement défini par

$$q_\tau(Y) = \inf \{y : F_Y(y) \geq \tau\}.$$

Si  $F_Y$  est continue, on retrouve la propriété intuitive<sup>2</sup>  $P(Y < q_\tau(Y)) = \tau$ . Les quantiles les plus couramment utilisés sont la médiane ( $\tau = 0,5$ ), les premier et dernier déciles ( $\tau = 0,1$  et  $\tau = 0,9$ ), et les premier et dernier quartiles ( $\tau = 0,25$  et  $\tau = 0,75$ ). Dans le langage courant, il y a parfois confusion entre la valeur du quantile et les personnes pour lesquelles la valeur de  $Y$  se situe en-dessous de ce quantile. Par exemple, on parle du premier décile pour désigner les 10 % de la population les moins riches. Cette désignation est incorrecte. Le premier décile désigne *stricto sensu* le seuil de revenus en-dessous duquel 10 % exactement de la population se situe.

### Modéliser l'ensemble de la distribution conditionnelle de la variable d'intérêt

Les régressions quantiles tentent alors d'évaluer comment les quantiles conditionnels  $q_\tau(Y|X)$ , définis par  $q_\tau(Y|X) = \inf \{y : F_{Y|X}(y) \geq \tau\}$ , se modifient lorsque les déterminants  $X = (1, X_2, \dots, X_p)'$  de  $Y$  varient. Il n'y a pas de raison en effet de supposer que l'impact d'une de ces caractéristiques  $X_k$  soit le même aux différents quantiles de la distribution conditionnelle de  $Y$ . On peut en trouver une illustration dans les classiques courbes de croissance utilisées dans les carnets de santé. Elles montrent comment la distribution du poids ou de la taille varie en fonction de l'âge. Plus précisément, elles représentent certains percentiles (traditionnellement les 3<sup>e</sup>, 25<sup>e</sup>, 75<sup>e</sup> et 97<sup>e</sup>) de ces distributions conditionnelles à l'âge. Elles permettent ainsi de vérifier que la croissance d'un enfant est normale en le situant dans la distribution correspondant à son âge. On observe sur de telles courbes que les différents percentiles des poids augmentent de manière à peu près linéaire avec l'âge sur les

premiers mois. Mais les taux de croissance correspondants diffèrent suivant le percentile : les droites ne sont pas parallèles.

Ce type de modélisation graphique est possible et utile lorsque l'on s'intéresse à un seul déterminant, mais atteint vite ses limites pour étudier simultanément l'effet de plusieurs caractéristiques sur la variable d'intérêt. Les régressions quantiles permettent justement d'étudier ce cadre multivarié : plus précisément, elles tentent de déterminer comment les quantiles de la distribution conditionnelle  $F_{Y|X}$  varient en fonction de  $X$ .

Dans la régression quantile standard, on suppose que ces quantiles de la distribution conditionnelle ont une forme linéaire :

$$q_\tau(Y|X) = X'\beta_\tau \quad (1)$$

où à chaque  $\tau$  correspond donc un vecteur de coefficients  $\beta_\tau = (\beta_{1,\tau}, \dots, \beta_{p,\tau})'$  correspondant aux  $p$  variables explicatives (dont la constante)<sup>3</sup>. Pour la suite, il peut être utile de remarquer que cette expression peut s'écrire de manière équivalente :

$$Y = X'\beta_\tau + \varepsilon_\tau, \text{ avec } q_\tau(\varepsilon_\tau | X) = 0 \quad (2)$$

La condition (1) est à rapprocher de celle qu'on postule dans la régression linéaire standard, dans laquelle on modélise la moyenne conditionnelle de la variable d'intérêt  $Y$  comme une expression linéaire des variables explicatives  $X$  :  $E(Y|X) = X'\beta$ . Une différence importante est qu'ici, on autorise les coefficients à différer d'un quantile à l'autre. Ceci apporte une information supplémentaire qui ne ressort pas d'une simple régression linéaire. Pour bien comprendre les implications de ce dernier point, considérons quelques exemples.

### Un premier cas particulier : le modèle de translation simple

Le premier exemple suppose que les variables explicatives n'ont d'impact que sur la moyenne

2. Cette égalité n'est cependant pas vraie dans le cas général, voir l'annexe pour une discussion de ce point et de quelques propriétés utilisées par la suite.

3. Cette dépendance linéaire n'exclut pas une dépendance plus compliquée des quantiles par rapport à certaines variables explicatives. Par exemple, dans l'exemple des courbes de croissance, les quantiles conditionnels à l'âge ne sont pas linéaires au-delà de la première année. En revanche, cette dépendance serait probablement bien approchée en utilisant non pas l'âge en niveau mais en logarithme, ou encore en utilisant une forme polynomiale de cette variable.

de la variable d'intérêt (et pas sur sa variance par exemple). Il s'agit du modèle de translation linéaire :

$$Y = X'\gamma + \varepsilon \quad (3)$$

où  $\varepsilon$  est indépendant de  $X$  et de moyenne nulle. Sous cette hypothèse, les résidus sont en particulier homoscedastiques, i.e.,  $V(\varepsilon | X) = \sigma^2$ . Dans ce modèle, les fonctions quantiles correspondant à différents  $\tau$  sont parallèles puisque  $q_\tau(Y | X) = X'\gamma + q_\tau(\varepsilon)$ . On est donc bien dans le cadre de l'hypothèse (1), mais ici seul le coefficient correspondant à la constante varie en fonction de  $\tau$  :  $\beta_{1,\tau} = \gamma_1 + q_\tau(\varepsilon)$ , tandis que  $\beta_{k,\tau} = \gamma_k$  pour  $k > 1$ . On parle dans ce cas d'*homogénéité des pentes*. On en trouvera une illustration dans le graphique I, dans le cas simple d'une régression univariée ( $X = (1, X_2)'$ ). Les droites correspondant aux régressions quantiles sont parallèles. Leur pente commune,  $\beta_{2,\tau} = \gamma_2$ , correspond à l'effet, constant ici, d'une augmentation de  $X_2$  de  $x_2$  à  $x_2 + 1$ , à  $\varepsilon$  fixé.

Cela signifie également que pour tout  $\tau$ , les coefficients  $(\beta_{k,\tau})_{k=2,\dots,p}$  sont les mêmes que ceux correspondant à la modélisation de la moyenne conditionnelle  $E(Y | X) = X'\gamma$ . Les résultats obtenus par régression quantile ou par une régression linéaire estiment donc les mêmes paramètres, par des méthodes différentes. Il peut néanmoins être intéressant d'utiliser les estimateurs des régressions quantiles plutôt que ceux des moindres carrés ordinaires. Ils sont plus robustes, par exemple, à la présence de valeurs aberrantes, et peuvent être plus précis pour certaines distributions de  $\varepsilon$ , point sur lequel nous reviendrons plus loin. Enfin, on ne sait pas en général si  $Y$  suit réellement un modèle de translation. Les résultats d'une régression quantile permettent de tester cette restriction. Elle implique en effet que les estimateurs de régressions quantiles effectuées pour différents  $\tau$  doivent être très proches.

### Un deuxième exemple : le modèle de translation-échelle

Le deuxième exemple, un peu plus général que le modèle de translation, suppose que les déterminants de la variable d'intérêt ont non seulement un impact sur sa moyenne mais aussi sur sa variance. Ces modèles, appelés translation-échelle, correspondent à une certaine forme d'hétéroscedasticité :

$$Y = X'\gamma + (X'\alpha)\varepsilon \quad (4)$$

avec encore une fois  $\varepsilon$  indépendant de  $X$ , de moyenne nulle et  $X'\alpha > 0$ . Dans un tel modèle, la dispersion de la variable dépendante est plus importante pour certaines valeurs de  $X$ . Un exemple classique est celui des salaires qui sont plus dispersés pour les diplômés du supérieur que pour les personnes sans diplôme (cf. le premier des deux exemples qui seront fournis dans la dernière section). Le modèle de translation-échelle correspondant à l'équation (4) implique que  $q_\tau(Y | X) = X'(\beta + q_\tau(\varepsilon)\alpha)$ . Ainsi, l'hypothèse (1) est bien vérifiée, avec  $\beta_\tau = \beta + q_\tau(\varepsilon)\alpha$ . L'impact des variables explicatives ne sera pas le même pour les différents quantiles, et il n'y a plus homogénéité des pentes. Dans l'exemple des salaires, l'effet du diplôme est faible pour les premiers quantiles ( $\beta_{k,\tau}$  petit pour  $\tau$  proche de 0) et plus fort pour les derniers quantiles ( $\beta_{k,\tau}$  grand pour  $\tau$  proche de 1).

Le modèle de translation-échelle est illustré par le graphique II, toujours dans le cas univarié. Ici les pentes  $\beta_{2,\tau}$  correspondant aux différentes régressions quantiles sont croissantes avec  $\tau$  (soit  $\alpha_2 > 0$ ), ce qui traduit une dispersion d'autant plus grande que  $X_2$  est élevé (comme dans l'exemple du diplôme). Cette information ne ressort pas d'une régression linéaire standard, qui se contente d'estimer le coefficient  $\gamma$  : on a en effet toujours  $E(Y | X) = X'\gamma$ .

### Un cadre plus général : le modèle à coefficients aléatoires\*

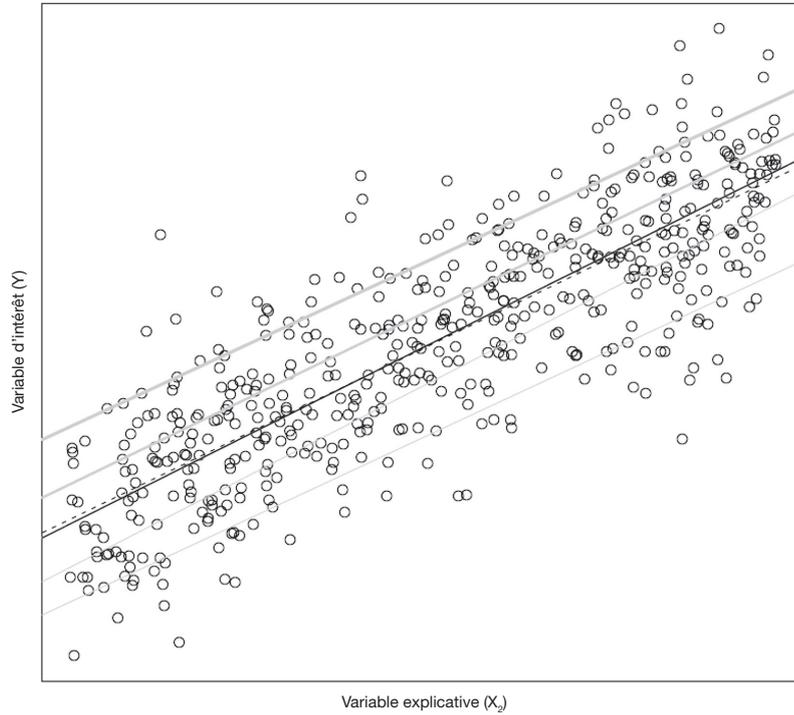
Une étape supplémentaire dans la généralisation est franchie par le modèle à coefficients aléatoires, qui s'écrit :

$$Y = X'\beta_U, \quad (5)$$

où  $U$  est indépendant de  $X$  et de loi uniforme sur  $[0, 1]$ . On suppose également que la fonction  $u \mapsto x'\beta_u$  est strictement croissante pour tout  $x^A$ . Dans ce modèle,  $U$  peut s'interpréter comme une composante individuelle inobservée qui positionne l'individu dans la distribution de  $Y$ . Par exemple, si on veut modéliser le niveau de salaire en fonction de caractéristiques observables comme le niveau d'étude,  $U$  pourrait correspondre à un niveau de productivité intrinsèque du salarié, qui implique en particulier que le niveau d'études ou d'autres caractéristiques ont des effets variant d'une personne à l'autre. Ce

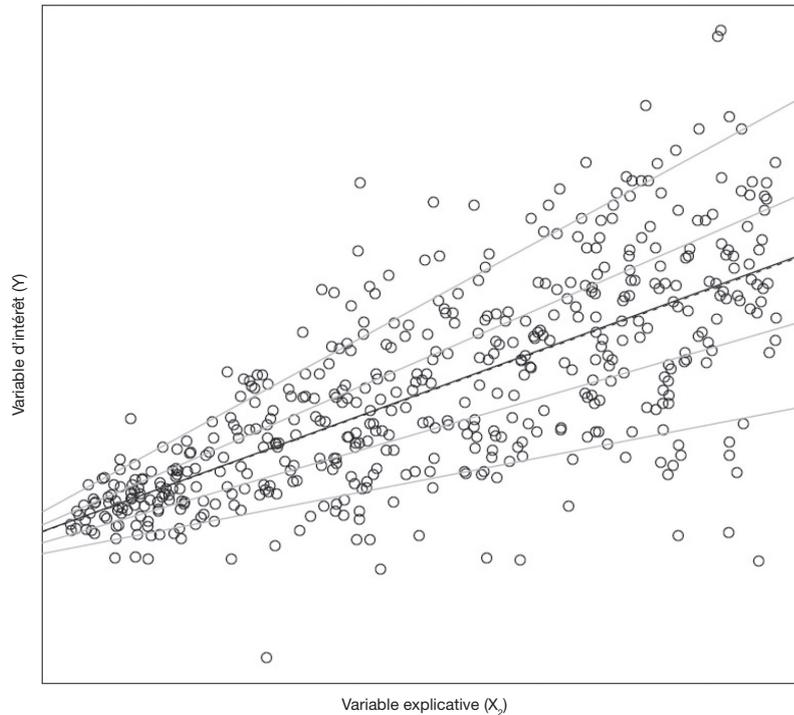
4. Ce modèle vérifie bien l'hypothèse (1) puisque, par indépendance de  $U$  et  $X$  et croissance de  $u \mapsto x'\beta_u$ ,  $P(Y \leq x'\beta_\tau | X = x) = P(x'\beta_U \leq x'\beta_\tau) = P(U \leq \tau) = \tau$ .

Graphique I  
Exemple de données distribuées selon un modèle de translation



Lecture: droites correspondant aux régressions quantiles pour les déciles d'ordre 1, 3, 7 et 9 (en gris), la médiane (en noir) et à une régression linéaire classique (en noir pointillé).  
Source : données simulées.

Graphique II  
Exemple de données distribuées selon un modèle de translation-échelle



Lecture : droites correspondant aux régressions quantiles pour les déciles d'ordre 1, 3, 7 et 9 (en gris), la médiane (en noir) et à une régression linéaire classique (en noir pointillé).  
Source : données simulées.

modèle à coefficients aléatoires repose sur des hypothèses très flexibles sur la dépendance en  $U$ . Il généralise les deux exemples précédents. Le modèle de translation linéaire correspond à un cas où les coefficients correspondant aux variables explicatives  $k > 1$ ,  $\beta_{k,U}$ , sont indépendants de  $U$ . Dans le modèle de translation-échelle, on a  $\beta_U = \gamma + q_U(\varepsilon)\alpha$ .

Ce modèle à coefficients aléatoires fournit une interprétation intéressante du coefficient  $\beta_\tau$  dans (1). Il implique en effet que si l'on modifie marginalement la variable observable  $X$  à  $U$  constant, l'effet sur la variable d'intérêt  $Y$  est égal à  $\beta_U$ . Ainsi,  $\beta_\tau$  correspond à l'effet marginal de  $X$  pour les individus au  $\tau^e$  quantile de la distribution des caractéristiques inobservées  $U$ . Si l'on s'intéresse par exemple à la modélisation des salaires en fonction du niveau d'études, la composante de  $\beta_{0,5}$  correspondant au nombre d'années d'études mesurera l'effet marginal d'une augmentation de celui-ci, pour les salariés dont la productivité intrinsèque  $U$  se situe à un niveau médian.

Cette interprétation n'est possible ici que parce qu'on a supposé l'existence d'un unique  $U$ , quelle que soit la valeur de  $X$ . Ceci implique la propriété dite d'invariance des rangs, qui signifie dans notre exemple que la productivité intrinsèque mentionnée plus haut est la même quel que soit le niveau d'études ou les autres variables observables utilisées. Il est important de comprendre que cette propriété restrictive, sur laquelle nous revenons plus loin, n'est pas requise dans les régressions quantiles. Toujours dans le cadre d'un modèle à coefficients aléatoires, on pourrait supposer par exemple que  $Y = X'\beta_{U_X}$ , avec  $U_x$  uniforme sur  $[0,1]$  pour tout  $x$  (la fonction  $u \mapsto x'\beta_u$  étant toujours supposée strictement croissante). Dans ce cas, l'hypothèse (1) serait encore vérifiée, mais plus la propriété d'invariance des rangs.

Le modèle à coefficients aléatoires, bien que plus général que les modèles de translation ou de translation-échelle, repose encore sur une forme très spécifique de dépendance entre  $Y$ ,  $X$  et l'effet inobservable. En toute généralité, cette dépendance peut s'écrire  $Y = g(X, U)$  où  $g$  est une fonction quelconque. Avec des hypothèses similaires à celle du modèle (5), le quantile conditionnel de  $Y$  correspond donc à  $q_\tau(Y|X) = g(X, \tau)$ . Le modèle de régression quantile (1) peut donc se voir comme une approximation linéaire de cette fonction  $g$ . Cette approximation est une des raisons pour laquelle on peut obtenir des estimations incohérentes.

En principe, la fonction  $\tau \mapsto g(x, \tau)$  doit être strictement croissante pour tout  $x$ , puisque  $\tau \mapsto q_\tau(Y|X=x)$  est strictement croissante. Mais la forme linéaire implique que si les pentes ne sont pas parallèles et si  $X$  a un large support, les différentes courbes  $x \mapsto x'\beta_\tau$  indicées par  $\tau$  se croiseront nécessairement pour certains quantiles  $\tau \neq \tau'^5$ .

### Interpréter les coefficients d'une régression quantile

Ces trois exemples aident à préciser l'apport de la régression quantile et ce que doit être la bonne lecture de ses résultats. La manière dont les distributions conditionnelles se modifient en fonction des variables explicatives renvoie à plusieurs questions qu'il importe de bien distinguer.

La première est simplement de décrire comment les quantiles conditionnels se modifient en fonction d'un changement d'un de ces déterminants. Les régressions quantiles sont l'outil adéquat pour répondre à cette question. Cette propriété suffit à en faire un outil précieux pour l'analyse des inégalités, ce qui a été un des premiers domaines d'application de la méthode (cf. encadré). Elle ne nécessite pas de supposer que ce sont des personnes comparables qui sont observées aux différents quantiles conditionnels. Prenons l'exemple de l'évaluation d'un programme spécifique d'apprentissage de la lecture. On cherche en particulier à comparer les scores à des épreuves standardisées des élèves dans des classes ayant bénéficié de ce programme avec ceux d'autres classes n'en ayant pas bénéficié. Si les deux groupes d'élèves étaient similaires initialement, ce qui est le cas par exemple si les classes ayant bénéficié du programme ont été tirées au sort, les régressions quantiles permettent de conclure sur l'impact de ce programme sur le premier décile. Autrement dit, on peut affirmer que grâce à ce programme, le niveau minimum atteint par 90 % des élèves a augmenté de  $x$  points.

À condition de faire l'hypothèse simple mais restrictive d'invariance des rangs, elles peuvent également permettre de répondre à une deuxième question, qui est de savoir comment  $Y$  varie suite à ce même changement d'un déterminant, pour

5. Même si les courbes théoriques ne se croisent pas, ceci peut également se produire sur les courbes estimées, en particulier lorsque l'échantillon sur lequel se base l'estimation est petit. Plusieurs solutions ont été proposées dans la littérature pour traiter de cette propriété indésirable (voir par exemple Chernozhukov et al., 2010).

les personnes se trouvant initialement à un certain niveau de la distribution conditionnelle de  $Y$ . Dans l'exemple du programme d'apprentissage de la lecture, il s'agit d'évaluer de combien augmente le score individuel des élèves dont le niveau en l'absence du programme expérimental se situait au niveau du premier décile. L'hypothèse d'invariance des rangs signifie ici qu'à l'intérieur d'une classe, les élèves sont toujours classés de la même manière entre eux, qu'ils bénéficient ou non du programme d'apprentissage. D'une manière générale, cette hypothèse n'est pas testable,

et il convient d'évaluer au cas par cas dans quelle mesure elle semble plausible pour pouvoir interpréter les résultats obtenus de cette manière.

Enfin, toujours sous l'hypothèse d'invariance des rangs, on peut également répondre à une troisième question, qui est celle de la distribution des effets pour chaque élève du programme d'apprentissage de la lecture sur le score individuel. Cette distribution permet de savoir, par exemple, que pour 10 % des élèves, le niveau du score a augmenté d'au moins  $x$  points.

Encadré

### RÉGRESSION QUANTILE ET ANALYSE DES INÉGALITÉS : QUELQUES PRÉCISIONS

Dans la mesure où elles permettent d'étudier l'ensemble de la distribution d'intérêt, les régressions quantiles constituent un outil privilégié pour l'étude des inégalités. Les estimations fournies par les régressions quantiles permettent par exemple de déterminer si les salaires sont plus ou moins dispersés parmi les détenteurs d'un niveau de diplôme par rapport à un autre niveau, en comparant les coefficients obtenus par des régressions quantiles pour les différents quantiles (voir par exemple l'illustration dans l'analyse des résultats de la régression quantile des salaires) : une mesure plus synthétique des inégalités intra groupes sera fournie par l'étude des différences interdéciles ( $q_{0,9}(Y|X) - q_{0,1}(Y|X)$ ) ou interquartiles ( $q_{0,75}(Y|X) - q_{0,25}(Y|X)$ ) que l'on peut reconstituer simplement grâce à l'approximation linéaire du quantile conditionnel. En évolution, on peut s'intéresser à l'évolution des déterminants des inégalités de salaires. L'un des articles pionniers sur cette question est Buchinsky (1994) qui s'intéresse à l'évolution des inégalités aux États-Unis sur longue période. Ce type d'analyse permet en particulier d'étudier comment les rendements de l'éducation et de l'expérience, les principaux facteurs explicatifs du revenu, ont évolué sur la période (par exemple du fait du progrès technique). Plus récemment, Charnoz *et al.* (2011) reproduisent ce type d'analyse pour la France. L'estimation à plusieurs dates permet par exemple d'étudier comment les inégalités de salaires ont évolué au cours du temps à l'intérieur d'un groupe de diplôme ou entre les groupes, en modélisant les différents quantiles prédits à partir des estimations obtenues, et à caractéristiques observables fixées. L'évolution temporelle des coefficients obtenus (correspondant aux rendements de l'éducation et de l'expérience dans une équation de Mincer classique) est en elle-même intéressante pour fournir des pistes d'interprétation des conséquences en termes d'inégalités des évolutions institutionnelles (salaire minimum en particulier) et économiques (progrès technique biaisé par exemple).

En revanche, il est important de rappeler que les régressions quantiles ne sont pas directement adaptées pour répondre à une question connexe, qui met

l'accent moins sur les facteurs économiques expliquant l'évolution des inégalités (comment les rendements de l'éducation évoluent avec le temps) que sur les effets de composition (les diplômés sont de plus en plus nombreux, or les salaires des diplômés sont plus dispersés). L'analyse de ces effets de composition a été popularisée sous le nom de décomposition d'Oaxaca-Blinder. En pratique, il s'agit de quantifier la part des différences observées dans les salaires de deux groupes (ici, deux dates différentes, mais ces méthodes sont souvent utilisées pour détecter d'éventuelles discriminations salariales, entre les hommes et les femmes par exemple) qui s'explique par des différences de composition. Cette décomposition repose sur le fait que pour deux groupes (noté  $A$  et  $B$ ), en supposant que les rendements des facteurs observables puissent être différents, on a :

$$E(Y_A) - E(Y_B) = (E(X_A) - E(X_B))\beta_A + E(X_B)(\beta_B - \beta_A)$$

où  $E(Z_i)$  correspond à la moyenne de la variable  $Z$  dans le groupe  $i$ . La première partie correspond à la part expliquée des inégalités (différences de structure entre les deux groupes), et le reste à la part inexpliquée (voir par exemple Meurs et Ponthieux (2006) ou Aeberhardt *et al.* (2010) pour une utilisation sur données françaises). En pratique, il s'agit donc d'estimer l'équation sur un des groupes, et d'utiliser les coefficients obtenus et les moyennes des caractéristiques observables dans les deux groupes pour obtenir la part expliquée.

Il n'est pas possible d'appliquer cette méthode à partir de régressions quantiles. La décomposition repose en effet sur la propriété de linéarité de l'espérance, que ne partage pas la fonction quantile :  $q_\tau(Y) \neq E(q_\tau(Y|X))$ . Dit autrement, le quantile de la variable d'intérêt sur l'ensemble de la population  $q_\tau(Y)$  ne correspond pas à la moyenne sur la population des quantiles conditionnels  $E(q_\tau(Y|X))$ . Sur les années récentes, de nombreuses solutions permettant d'étendre les méthodes de décomposition aux différents quantiles ont été proposées. Leur présentation dépasse le cadre de cet article. On en trouvera une synthèse dans Fortin *et al.* (2011).

Pour bien comprendre que les réponses à ces trois questions peuvent être différentes, il peut être utile de considérer une situation très simple (cf. tableau). Supposons que la population d'élèves est composée de cinq types en proportions identiques, qu'on identifiera par les lettres de A à E. Par ailleurs,  $X = (1, X_2)$ , où  $X_2 = 1$  si l'élève a bénéficié du programme spécifique de lecture, 0 sinon.  $X_2$  est supposé tiré au sort, si bien qu'on observe les différents types dans les mêmes proportions lorsque  $X_2 = 0$  et  $X_2 = 1$ . Suivant la valeur de  $X_2$ , les élèves peuvent avoir des valeurs différentes de  $Y$ . Les effets du programme d'apprentissage correspondent pour un élève à un passage de  $X_2 = 0$  à  $X_2 = 1$ . Une régression quantile d'ordre 0,5 (*régression médiane*) de  $Y$  sur  $X_2$  mesure l'écart entre la médiane de la distribution de  $Y$  conditionnelle à  $X_2 = 0$  et la médiane de la distribution conditionnelle à  $X_2 = 1$ . Il vaut donc ici 2 (6 - 4). Cette valeur est différente de l'effet de passer de  $X_2 = 0$  à  $X_2 = 1$  pour les élèves qui sont dans le groupe médian quand  $X_2 = 0$ . Ces élèves sont ceux de type C pour qui  $\Delta Y^C = 1$ .

Cette différence vient du fait qu'ici, les élèves ne sont pas ordonnés (en termes de  $Y$ ) de la même manière lorsque  $X_2 = 0$  et  $X_2 = 1$ , autrement dit que l'hypothèse d'invariance des rangs n'est pas vérifiée. Dans notre exemple artificiel, certains élèves sont plus à l'aise avec la méthode spécifique d'apprentissage de la lecture, d'autres moins. Ainsi, les élèves de type C sont devant ceux de type B lorsque  $X_2 = 0$  mais derrière eux lorsque  $X_2 = 1$ . Ceci interdit d'interpréter le coefficient de la régression quantile d'ordre 0,5 comme l'effet d'un passage de  $X_2 = 0$  à  $X_2 = 1$  pour les élèves qui sont dans le groupe médian quand  $X_2 = 0$ . Enfin, l'effet individuel du changement

de méthode,  $\Delta Y$ , a une médiane de 3. On constate ici que la médiane de la distribution de  $\Delta Y$  ne correspond pas à la différence des médianes des distributions de la variable d'intérêt, ou, en termes mathématiques,  $q_{0,5}(\Delta Y) \neq q_{0,5}(Y | X_2 = 1) - q_{0,5}(Y | X_2 = 0)$ .

Ces remarques permettent de bien cadrer les usages qui peuvent être faits, ou non, des résultats d'une régression quantile. Il est utile de faire le lien avec ceux obtenus par une régression linéaire classique. L'objet principal de celle-ci est de modéliser la manière dont la moyenne conditionnelle varie en fonction de déterminants : sur notre exemple, cela correspondrait à 3 (= 7,2 - 4,2). Du fait de la linéarité de la moyenne, la différence des moyennes conditionnelles correspond également à l'effet moyen de l'augmentation de  $X_2$ , c'est-à-dire à la moyenne de  $\Delta Y$ . En revanche, cette différence ne correspond pas à ce que gagnerait un élève, dont la valeur de  $Y$  est proche de la moyenne conditionnelle lorsque  $X_2 = 0$  (dans notre exemple, il s'agit encore de C), s'il passait à  $X_2 = 1$ . Contrairement aux régressions quantiles, la condition d'invariance des rangs n'est pas suffisante pour permettre une telle interprétation.

### Des estimateurs également plus adaptés à certains types de variables

Indépendamment de sa capacité à aller « au-delà » de la moyenne, la régression quantile peut aussi être préférée à la régression linéaire au nom d'arguments plus techniques, spécifiques à certains types de variables ou de modèles.

Une première raison est que la régression quantile est robuste aux valeurs aberrantes de la

Tableau  
Exemple d'effets d'une variable dichotomique  $X_2$  sur la distribution d'une variable d'intérêt  $Y$ , pour une population fictive divisée en cinq types.

Type	Valeur de $Y$ si $X_2 = 0$	Valeur de $Y$ si $X_2 = 1$	$\Delta Y$
	(1)	(2)	(2) - (1)
A	1	4	3
B	2	6	4
C	4	5	1
D	5	11	6
E	9	10	1
Médiane	4	6	3
Moyenne	4,2	7,2	3

Lecture : les cinq types sont supposés en proportion égale dans la population. Conditionnellement au fait d'avoir  $X_2 = 0$  (respectivement  $X_2 = 1$ ), la valeur médiane de  $Y$  est égale à 4 (respectivement 6). Le résultat d'une régression quantile de l'effet de  $X_2$  sur la médiane serait ainsi de 6 - 4 = 2. Ce résultat est donc différent de la médiane de l'effet de  $X_2$  sur  $Y$  ( $\Delta Y$ ), qui vaut 3.

variable d'intérêt ou à la présence d'erreurs très dispersées. Intuitivement, cette propriété est due au fait que les quantiles sont moins sensibles que la moyenne à la présence de valeurs très grandes<sup>6</sup>. Supposons que la variable  $Y^*$  vérifie le modèle de translation simple (3), mais que dans de très rares cas, les données observées ne correspondent pas à la variable d'intérêt  $Y^*$  mais à une variable erronée, éventuellement corrélée avec les variables explicatives  $X$ . Formellement, on observe donc  $Y = AX'\delta + (1-A)Y^*$ , où  $A$  est une variable inobservée valant 1 lorsque  $Y$  est aberrant, 0 sinon, avec  $P(A=1 | X, \varepsilon) = p$  petit. Une régression linéaire de la variable observée (avec erreur)  $Y$  sur  $X$  donnera une estimation biaisée de notre paramètre d'intérêt  $\gamma$ , puisqu'elle sera égale, pour un échantillon de taille tendant vers l'infini, à  $\gamma_{MCO} = \gamma + p(\delta - \gamma)$ . Si  $\delta$  est très différent de  $\gamma$ , le terme de biais peut être important même lorsque la probabilité d'observer des valeurs erronées  $p$  est faible. En revanche, on montre en annexe que si les valeurs aberrantes  $X'\delta$  sont très grandes, l'estimateur de l'effet de  $X_k$  ( $k > 1$ ) obtenu par une régression quantile vaut bien  $\gamma_k$ . En d'autres termes, la présence de valeurs aberrantes n'affecte pas les résultats de la régression quantile, sauf les coefficients de la constante. Dans le cas plus général, si on suppose que  $Y^*$  vérifie (5), cette propriété n'est plus exactement vérifiée, mais elle reste une bonne approximation dès lors que la proportion  $p$  de valeur aberrantes est faible. On obtient en effet par la régression quantile le coefficient  $\beta_{\tau/(1-p)}$  au lieu de  $\beta_{\tau}$  (voir l'annexe).

Par ailleurs, même lorsque les données sont toutes correctement observées, la distribution sous-jacente peut être telle qu'une régression linéaire n'est pas adaptée. C'est le cas en particulier lorsque  $\varepsilon$ , dans le modèle de translation simple, peut prendre des valeurs très grandes avec une probabilité non négligeable (on parle de distributions à queue épaisse). Ainsi, lorsque l'on mesure les patrimoines ou les très hauts revenus, les résidus peuvent être très dispersés au sens où l'on peut observer des valeurs très élevées, qui ne sont pas des valeurs aberrantes. C'est la raison pour laquelle on les modélise en général par des lois de Cauchy ou certaines lois de Pareto, qui n'ont pas d'espérance. Dans ce cas, l'estimateur des moindres carrés ordinaires n'est pas convergent : même pour des échantillons énormes, il pourra prendre des valeurs très différentes du vrai paramètre  $\beta$ . À l'inverse, l'estimateur obtenu par régression quantile sera convergent car intuitivement, cet estimateur ne dépend pas des queues de distribution de  $\varepsilon$  mais

seulement de la distribution de  $\varepsilon$  autour du quantile considéré (cf. *infra*).

Enfin, outre leur robustesse, une propriété importante des quantiles est qu'ils sont invariants par une transformation monotone : si  $g$  est une fonction croissante continue à gauche, on a  $q_{\tau}(g(Y)) = g(q_{\tau}(Y))$  (une démonstration est donnée en annexe). Cette propriété n'est, bien sûr, pas vérifiée par l'espérance. Comme on le présentera en détail plus loin, ceci rend les régressions quantiles plus naturelles et simples à utiliser dans des modèles non-linéaires<sup>8</sup> comme les modèles à censure fixe, ou les modèles de durées (voir par exemple Biliias et Koenker, 2001 ou Fitzenberger et Wilke, 2005). On peut trouver aussi des applications aux modèles de comptage (Machado et Silva, 2005).

## Techniques d'estimation et principales propriétés statistiques\*

Ces différents apports de la régression quantile ayant été précisés, on peut fournir quelques détails sur la façon dont elle est estimée en pratique et les principales propriétés des estimateurs obtenus, notamment en termes de précision.

### Définition de l'estimateur et propriétés statistiques

Pour bien comprendre le principe des régressions quantiles, il est utile de détailler comment on peut estimer les quantiles d'une variable d'intérêt  $Y$  à partir d'un échantillon  $(Y_i)_{i=1 \dots n}$  de variables supposées i.i.d. La manière la plus intuitive de calculer l'estimateur standard  $\hat{q}_{\tau}(Y)$  consiste à ordonner ces  $n$  variables, le quantile d'ordre  $\tau$  étant fourni par la  $[n\tau]^e$  observation où  $[n\tau]$  est le plus petit entier supérieur ou égal à

6. Par exemple, les deux échantillons (1,5,6,8,11) et (1,5,6,8,999999) ont la même médiane mais des moyennes très différentes.

7. On néglige ici les erreurs d'estimation liées au fait que l'échantillon est de taille finie. Autrement dit, nos résultats doivent se comprendre comme étant les valeurs limites des estimateurs lorsque la taille de l'échantillon tend vers l'infini.

8. Cette propriété signifie aussi qu'on pourra facilement déduire l'effet marginal d'un déterminant du salaire (par exemple) sur un de ces quantiles conditionnels à partir d'une régression quantile modélisant le logarithme du salaire. On aura en effet  $q_{\tau}(W|X) = \exp(X\beta_{\tau})$ , où  $\beta_{\tau}$  correspond au coefficient estimé par la régression quantile de  $\log(W)$  sur  $X$ .

$n\tau^9$ . Mais il est plus utile, pour le passage aux régressions quantiles, de remarquer qu'on a également<sup>10</sup> (cf. démonstration donnée en annexe) :

$$\hat{q}_\tau(Y) = \arg \min_b \frac{1}{n} \sum_{i=1}^n \rho_\tau(Y_i - b). \quad (6)$$

où  $\rho_\tau(\cdot)$  est une fonction test définie par  $\rho_\tau(u) = (\tau - 1 \{u < 0\})u$ . Par exemple, pour  $\tau = 1/2$ , c'est-à-dire si l'on s'intéresse à la médiane, la fonction test correspond simplement à la (demi-) valeur absolue.

L'intérêt de cette définition est de s'étendre simplement au cadre conditionnel qui nous intéresse. Il suffit en effet de remplacer  $\hat{q}_\tau(Y)$  et  $b$  dans (6) par respectivement  $q_\tau(Y|X)$  et une fonction  $b(X)$ . Dans le cas des régressions quantiles classiques, on peut se limiter aux fonctions linéaires puisqu'on suppose que  $q_\tau(Y|X) = X'\beta_\tau$ . On a alors :

$$\beta_\tau = \arg \min_\beta E[\rho_\tau(Y - X'\beta)]. \quad (7)$$

On peut noter l'analogie avec le modèle de régression linéaire classique, qui modélise l'espérance conditionnelle de  $Y$  par une forme linéaire en  $X$  :  $E(Y|X) = X'\beta_0$ . L'espérance d'une variable aléatoire pouvant être obtenue par  $E(Y) = \arg \min_a E[(Y - a)^2]$ , le coefficient  $\beta_0$  est défini par  $\beta_0 = \arg \min_\beta E[(Y - X'\beta)^2]$ . La fonction de perte quadratique qui est utilisée dans une régression linéaire par les moindres carrés ordinaires est donc remplacée, dans la régression quantile, par la fonction test  $\rho_\tau(\cdot)$ . Celle-ci augmentant de manière linéaire et non quadratique avec le résidu, les très grands écarts sont beaucoup moins pénalisés, ce qui explique la robustesse de la régression quantile aux valeurs extrêmes ou aberrantes discutée plus haut.

L'estimateur de la régression quantile est alors obtenu en remplaçant l'espérance dans (7) par la moyenne sur l'échantillon :

$$\hat{\beta}_\tau = \arg \min_\beta \frac{1}{n} \sum_{i=1}^n \rho_\tau(Y_i - X_i'\beta). \quad (8)$$

Dans le cas où  $\tau = 1/2$ , l'estimateur revient à minimiser la somme des valeurs absolues des erreurs  $Y_i - X_i'\beta$ . On parle de l'estimateur des moindres déviations absolues (*Least Absolute Deviation Estimator*), qui est de fait utilisé depuis longtemps comme une alternative robuste aux moindres carrés ordinaires.

On insiste sur le fait que cette estimation se fait en utilisant l'ensemble de l'échantillon.

Contrairement à une confusion fréquente, elle ne consiste pas à diviser l'échantillon en fonction des quantiles de la variable d'intérêt, puis à effectuer des régressions linéaires sur les sous-échantillons ainsi obtenus. Procéder à des régressions séparées sur ces sous-groupes ne serait pas cohérent puisque ceci reviendrait à contraindre les valeurs inférieure et supérieure de la variable d'intérêt au sein de chaque groupe, puis à étudier comment varie cette même variable d'intérêt en fonction de ses variables explicatives. Cette confusion en recoupe une autre : celle qui est souvent faite entre le niveau des quantiles (les bornes de l'intervalle) et les personnes dont la valeur de la variable d'intérêt se situe dans ces intervalles.

Cette estimation appliquée à l'ensemble de l'échantillon peut alors se faire pour tout quantile d'ordre  $\tau$ , où  $\tau \in [0,1]$ . Il existe donc en principe une infinité de régressions quantiles possibles. En pratique, le nombre de quantiles qu'on estime dépendra de la taille de l'échantillon. Il est bien entendu illusoire de tenter d'approcher très finement une distribution avec un nombre fini d'observations : le nombre de quantiles empiriques distincts est restreint (de l'ordre de  $n \ln(n)$  où  $n$  désigne la taille de l'échantillon, voir Portnoy (1992)). Le choix de modéliser l'ensemble des percentiles ou simplement les quartiles et la médiane dépendra non seulement du degré de précision souhaitée pour décrire la distribution, mais aussi des données disponibles.

Enfin, un indicateur de la qualité de l'ajustement d'une régression quantile a été proposé par Koenker et Machado (1999). Il est défini par

$$R^1(\tau) = 1 - \frac{\min_{b \in \mathbb{R}^p} \rho_\tau(Y_i - X_i'b)}{\min_{b_0 \in \mathbb{R}} \rho_\tau(Y_i - b_0)}.$$

Comme le  $R^2$ , ce critère est compris entre 0, lorsque l'estimateur des coefficients relatifs à  $(X_2, \dots, X_p)$  vaut 0, et 1, lorsque  $Y$  est une fonction linéaire déterministe de  $X$ . Il augmente également lorsqu'on ajoute des variables explicatives au modèle.

9. Cet estimateur n'est pas le seul utilisé. Le logiciel R propose ainsi rien moins que neuf estimateurs différents, basés sur les définitions données dans Hyndman et Fan (1996).

10. En toute rigueur, il n'y a pas toujours unicité au programme de minimisation  $\min_a E[r_\tau(Y-a)]$  (cf. l'annexe pour une discussion). On néglige ici ces complications.

## Algorithmes utilisés

Il n'existe pas de solution explicite à (8), si bien qu'il faut résoudre ce programme numériquement. Or la fonction objectif n'est pas différentiable, puisque la fonction  $\rho_\tau$  n'est pas dérivable en 0. Les algorithmes standards tels que celui de Newton Raphson ne peuvent donc pas être utilisés directement. Cependant, (8) peut se réécrire comme un programme linéaire :

$$\min_{(\beta, u, v) \in \mathbb{R}^p \times \mathbb{R}_+^{2n}} \tau \mathbf{1}'u + (1-\tau) \mathbf{1}'v \text{ s.t. } X\beta + u - v - Y = 0,$$

où  $X = (X_1, \dots, X_n)'$ ,  $Y = (Y_1, \dots, Y_n)'$  et  $\mathbf{1}$  est un vecteur de 1 de taille  $n$ . Les variables supplémentaires  $u$  et  $v$  désignent respectivement les parties positives et négatives des résidus. Cette reformulation, a priori peu naturelle car elle augmente la dimension du vecteur à optimiser, est très utile car de nombreux algorithmes permettent de résoudre de tels programmes linéaires. La méthode du simplexe a été jusqu'à récemment la méthode la plus classique pour résoudre ce type de problèmes linéaires. Cependant, elle devient coûteuse en temps de calcul lorsque le nombre d'observations augmente et elle n'est donc adaptée qu'à de petits échantillons. Pour des échantillons plus conséquents, les méthodes de points intérieurs sont plus performantes pour résoudre ces problèmes (Portnoy et Koenker, 1997). On trouvera dans Givord et D'Haultfœuille (2013) une description pratique de l'implémentation de ces méthodes dans les logiciels statistiques standards.

## Propriétés asymptotiques de l'estimateur et estimation de la précision

Les propriétés asymptotiques de  $\hat{\beta}_\tau$  sont délicates à établir car, contrairement à l'estimateur des moindres carrés, il n'existe pas de forme explicite pour  $\hat{\beta}_\tau$ . Pour plus de détails, on se référera par exemple à l'ouvrage de Koenker (2005). Nous nous limitons ici au résultat principal sur la loi asymptotique de  $\hat{\beta}_\tau$ .

**Théorème :** *Supposons que  $\varepsilon_\tau = Y - X'\beta_\tau$  admette, conditionnellement à  $X$ , une densité en 0  $f_{\varepsilon_\tau|X}(0|X)$  et que  $J_\tau = E\left[f_{\varepsilon_\tau|X}(0|X)XX'\right]$  soit inversible. Alors*

$$\sqrt{n}\left(\hat{\beta}_\tau - \beta_\tau\right) \xrightarrow{d} \mathbf{N}\left(0, \tau(1-\tau)J_\tau^{-1}E[XX']J_\tau^{-1}\right) \quad (9)$$

Dans le cas du modèle de translation simple (3), la variance asymptotique prend une forme

particulièrement simple. On a en effet,  $\varepsilon_\tau = \varepsilon - q_\tau(\varepsilon)$  et la variance asymptotique  $V_{as}$  s'écrit plus simplement

$$V_{as} = \frac{\tau(1-\tau)}{f_\varepsilon(q_\tau(\varepsilon))^2} E[XX']^{-1}.$$

Cette variance est très proche de celle des MCO avec résidus homoscédastiques, si ce n'est que  $\sigma^2 = V(\varepsilon)$  est remplacé par  $\tau(1-\tau)/f_\varepsilon(q_\tau(\varepsilon))^2$ . Le terme de densité est logique : seuls les résidus autour de  $q_\tau(\varepsilon)$  vont apporter de l'information sur la valeur du quantile conditionnel de  $Y$ . Ce résultat explique que, même dans le cas de translation simple, il peut être parfois préférable d'utiliser une régression quantile pour certaines distributions des termes inobservés  $\varepsilon$ . L'estimation par régression quantile sera plus précise qu'une estimation par MCO lorsque  $\tau(1-\tau)/f_\varepsilon(q_\tau(\varepsilon))^2 < \sigma^2$ .

En dehors du modèle restrictif de translation simple, la variance asymptotique est plus complexe à estimer que dans le cadre d'un modèle de régression linéaire simple. Plusieurs méthodes d'inférence ont été proposées pour construire des tests ou des intervalles de confiance sur  $\beta_\tau$ , et il n'existe pas à l'heure actuelle de consensus sur la méthode à utiliser. On trouvera dans Kocherginsky *et al.* (2005) une présentation générale de ces méthodes et une discussion pratique des cas où certaines sont plus ou moins indiquées. Le choix dépend des hypothèses plus ou moins restrictives qu'on accepte de faire sur le modèle sous-jacent (modèle de translation...), de la taille de l'échantillon ou du nombre de variables du modèle.

Certaines méthodes s'appuient sur une estimation directe de la variance asymptotique en partant de la formule (9). La difficulté principale de cette approche est la présence de la densité conditionnelle  $f_{\varepsilon_\tau|X}(0|X)$ , qui est délicate à estimer (cf. annexe pour plus de détails). Dans le cadre restrictif d'un modèle de translation-échelle, une méthode basée sur les tests de rang est parfois utilisée (cf. Koenker, 2005). Il est surtout courant de s'appuyer sur des méthodes de *bootstrap*. Elles consistent à générer des échantillons factices par des tirages avec remise à partir de l'échantillon initial, et à effectuer une régression quantile sur ces échantillons (cf. annexe). L'inconvénient de ces méthodes est qu'elles sont souvent coûteuses en temps de calcul. Ce dernier augmente à la fois avec la taille de l'échantillon et le nombre de variables explicatives. Une solution récente (Markov

*Chain Marginal Bootstrap*, ou MCMB) a été proposée par He et Hu (2002) pour résoudre en partie ce problème quand le nombre de variables explicatives est important.

Enfin, l'un des intérêts de la régression quantile étant de ne pas supposer *a priori* que les variables explicatives ont un effet homogène sur l'ensemble de la distribution de la variable d'intérêt, il est tout à fait possible de tester cette hypothèse à partir des estimations obtenues. Par exemple, l'homogénéité de l'effet de l'une des variables  $X_k$  correspond à l'égalité des coefficients  $\beta_{k,\tau_1}, \dots, \beta_{k,\tau_m}$  (où  $(\tau_1, \dots, \tau_m)$  peuvent être par exemple l'ensemble des déciles), ce qui peut se tester simplement. Un tel test s'appuie sur la distribution jointe asymptotique de  $(\hat{\beta}_{\tau_1}, \dots, \hat{\beta}_{\tau_m})$ , donnée par le résultat suivant :

$$\sqrt{n} \left( \hat{\beta}_{\tau_k} - \beta_{\tau_k} \right)_{k=1}^m \xrightarrow{d} \mathbf{N}(0, V), \quad (10)$$

où  $V$  est une matrice par bloc dont le bloc  $V_{k,l}$  vérifie

$$V_{k,l} = [\min(\tau_k, \tau_l) - \tau_k \tau_l] J_{\tau_k}^{-1} E[XX'] J_{\tau_l}^{-1}.$$

Ce résultat est donc une généralisation du théorème précédent à plusieurs quantiles.

## Plusieurs extensions\*

Le lien entre régression ordinaire et régression quantile conduit à retrouver, dans le contexte de la régression quantile, un certain nombre de questions économétriques classiques : l'endogénéité des variables explicatives, la mobilisation de données de panel, l'application à l'évaluation des politiques publiques, les non-linéarités... Certaines de ces questions sont assez complexes et relèvent plutôt, à ce stade, du domaine de la recherche. On va néanmoins tenter d'en donner un aperçu rapide.

### Les régressions quantiles instrumentales

Comme en régression linéaire, il arrive fréquemment que certaines composantes des variables  $X$  soient *a priori* endogènes. Par exemple, dans une étude sur l'impact d'un dispositif de formation sur le salaire, le fait de participer à ce dispositif peut être lié à des caractéristiques inobservées qui influent

également sur le salaire. Dans ce cas, l'estimateur  $\hat{\beta}_\tau$  défini par (8) ne mesure pas l'effet causal du dispositif de formation.

En revanche, on peut disposer d'instruments affectant ces variables mais pas directement les composantes inobservées de la variable d'intérêt (représentées par le résidu  $\varepsilon_\tau = \varepsilon - q_\tau(\varepsilon)$ ). Plus précisément, on se place dans le cadre de la régression quantile classique, définie par (2),

$$Y = X'\beta_\tau + \varepsilon_\tau,$$

mais on suppose que certaines des variables explicatives, notées  $X_1 \in \mathbb{R}^q$ , sont endogènes, c'est-à-dire telles que  $q_\tau(\varepsilon_\tau | X_1) \neq 0$ . Les autres variables explicatives, notées  $X_2$ , sont supposées exogènes (c'est-à-dire qu'elles vérifient  $q_\tau(\varepsilon_\tau | X_2) = 0$ ). On suppose disposer par ailleurs d'instruments, notés  $Z_1 \in \mathbb{R}^r$  (avec  $r \geq q$ ).  $Z_1$  doit être corrélé aux variables explicatives endogènes mais pas aux résidus, si bien que :

$$q_\tau(\varepsilon_\tau | Z) = 0 \quad (11)$$

avec  $Z = (Z_1', X_2')'$ . Cette hypothèse est l'équivalent de l'hypothèse  $E(\varepsilon | Z) = 0$  en régression linéaire instrumentale. Mais il serait tout à fait incorrect de transposer la méthode classique des doubles moindres carrés, qui consiste à régresser la variable d'intérêt non sur la variable endogène, mais sur sa projection sur les variables instrumentales obtenues par une première étape. Cette méthode repose en effet sur la propriété de linéarité de l'espérance, que ne vérifient pas les quantiles. En revanche, on peut proposer une méthode s'appuyant sur la condition d'exclusion (11). En effet, cette condition implique que :

$$q_\tau(Y - X_1'\beta_{1\tau} | Z) = q_\tau(X_2'\beta_{2\tau} + \varepsilon_\tau | Z_1, X_2) \quad (12) \\ = X_2'\beta_{2\tau} + Z_1' 0_r$$

où  $0_r$  est un vecteur de zéros de taille  $r$ . L'équation (12) est à la base d'une méthode proposée par Chernozhukov et Hansen (2008). Elle signifie que dans une régression quantile de  $Y - X_1'\beta_{1\tau}$  sur  $X_2$  et  $Z_1$ , les coefficients de  $Z_1$  sont égaux à 0. Par ailleurs, si l'on se trompe en considérant un coefficient  $\beta \neq \beta_{1\tau}$ , la variable modifiée  $Y - X_1'\beta$  sera toujours liée à  $X_1$ . Donc, si  $X_1$  et  $Z_1$  sont dépendants, ce qui correspond à la condition de rang dans les régressions linéaires, la régression quantile de  $Y - X_1'\beta$  sur  $Z$  conduira *a priori* à un coefficient non nul sur

$Z_1$ <sup>11</sup>. L'idée de Chernozhukov et Hansen est alors d'inverser la régression quantile, en estimant  $\beta_{1\tau}$  par le paramètre  $\hat{\beta}_{1\tau}$  permettant d'obtenir, dans la régression quantile de  $Y - X_1'\hat{\beta}_{1\tau}$  sur  $Z$ , un coefficient égal à 0 pour  $Z_1$ . En pratique, les auteurs proposent l'algorithme suivant:

1. Définir une grille sur  $\beta_{1\tau}$ ,  $\{b_1, \dots, b_J\}$ .
2. Pour  $j = 1$  à  $J$ :
  - Calculer les estimateurs de régression quantile de  $Y - X_1' b_j$  sur  $(Z_1, X_2)$ . Soient  $(\hat{\gamma}(b_j), \hat{\beta}_{2\tau}(b_j))$  les estimateurs correspondants.
  - Calculer la statistique de Wald correspondant au test de  $\gamma(b_j) = 0$ :

$$W_n(b_j) = n\hat{\gamma}(b_j)' \hat{V}_{as}^{-1}(\hat{\gamma}(b_j)) \hat{\gamma}(b_j),$$

où  $\hat{V}_{as}^{-1}(\hat{\gamma}(b_j))$  est l'estimateur de la variance asymptotique de  $\hat{\gamma}(b_j)$ .

3. Définir l'estimateur de  $\beta_\tau = (\beta_{1\tau}, \beta_{2\tau})$  par :

$$\hat{\beta}_{1\tau} = \arg \min_{j=1 \dots J} W_n(b_j), \hat{\beta}_{2\tau} = \hat{\beta}_{2\tau}(\hat{\beta}_{1\tau}).$$

La commande Stata *ivqreg*, introduite récemment<sup>12</sup>, utilise cet algorithme pour estimer  $\beta_\tau$  dans un tel modèle. Même en l'absence de procédure préprogrammée, cet algorithme a l'intérêt de ne s'appuyer que sur des régressions quantiles classiques. Il peut donc être mis en œuvre simplement avec des logiciels standards. En pratique, la grille doit être suffisamment fine pour ne pas altérer les propriétés asymptotiques de l'estimateur (pour plus de détails, se reporter à Chernozhukov et Hansen (2008)). Pour que le temps de calcul reste raisonnable, le nombre de variables endogènes doit donc être petit ( $q = 1$  ou  $2$ ).

Notons enfin que d'autres solutions existent pour estimer des régressions quantiles instrumentales. Abadie *et al.* (2002), en particulier, proposent de recourir à une approche par régression quantile pondérée pour les cas où la variable endogène  $X_1$  et l'instrument  $Z_1$  sont binaires. Quelle que soit la méthode retenue, la difficulté principale est évidemment de trouver un instrument valide, c'est-à-dire vérifiant (11). Nous proposons, dans la dernière section de l'article, un exemple d'instrument utilisant les données issues d'une expérimentation sociale tiré d'Abadie *et al.* (2002).

## Les régressions quantiles avec des données de panel

L'utilisation de données de panel, c'est-à-dire de données répétées pour les mêmes unités, peut être également une manière de traiter la présence d'hétérogénéité individuelle inobservée. Ces données sont en effet plus souvent disponibles que des variables instrumentales valides. Sous l'hypothèse, certes restrictive, que l'hétérogénéité individuelle est fixe dans le temps et indépendante des termes résiduels, une simple différenciation permet dans le cadre des régressions linéaires de se débarrasser de ces effets individuels fixes dans le temps. Cependant, ces méthodes ne s'appliquent plus directement dans le cas des quantiles. Ces derniers n'ont pas de propriété de linéarité comme la moyenne. Les quantiles des variables différenciées ne correspondent pas directement aux quantiles d'intérêt.

Sur la période très récente, de nombreux estimateurs permettant de tenir compte de la présence d'effets individuels dans des régressions quantiles ont été proposés. L'estimateur de Canay (2011) a l'avantage de s'appuyer sur des techniques standards utilisant des procédures couramment disponibles dans les logiciels statistiques. Pour en comprendre le principe, il est utile d'utiliser les notations du modèle à coefficients aléatoires présenté plus haut, dans lequel on prend en compte également un effet individuel  $\alpha_i$  fixe dans le temps :

$$Y_{it} = X_{it}'\beta_{U_{it}} + \alpha_i, \quad (13)$$

où  $\alpha_i$  et  $U_{it}$  sont inobservables et  $U_{it}$  est indépendant de  $(\alpha_i, X_{it})$  et suit une distribution uniforme sur  $[0,1]$ .  $X_{it}$  représentent les variables explicatives tandis que  $\alpha_i$  représente des caractéristiques fixes dans le temps. Notons qu'aucune hypothèse n'est imposée ici quant à la dépendance entre  $X_{it}$  et  $\alpha_i$ . Ce modèle permet donc de résoudre en partie le problème d'endogénéité puisqu'il autorise des corrélations entre  $X_{it}$  et les facteurs inobservables individuels, pourvu que ces derniers soient stables dans le temps.

11. La condition de dépendance exacte nécessaire entre  $X_i$  et  $Z$  est plus difficile à expliciter que dans les modèles linéaires, voir Chernozhukov et Hansen (2008) pour plus de détails à ce sujet.

12. Package développé par D.W. Kwak, disponible à l'adresse <http://faculty.chicagobooth.edu/christian.hansen/research/>.

En introduisant  $e_{it} = X'_{it}(\beta_{U_{it}} - \beta_\tau)$ , on peut réécrire le modèle comme :

$$Y_{it} = X'_{it}\beta_\tau + \alpha_i + e_{it}, \text{ avec } q_\tau(e_{it} | X_i) = 0. \quad (14)$$

Sous des hypothèses techniques détaillées par Canay (2011), on peut montrer que les termes  $\beta_\tau$  sont identifiés et proposer un estimateur convergent lorsque  $T$  tend vers l'infini. Plus précisément, Canay propose un estimateur basé sur les deux étapes suivantes :

1. Estimation par un estimateur *within* classique de la régression linéaire

$$Y_{it} = X'_{it}\beta_\mu + \alpha_i + u_{it},$$

avec :

$$E(u_{it} | X_{it}, \alpha_i) = 0 \text{ et } \beta_\mu = E[\beta_{U_{it}}].$$

À partir de l'estimation de  $\beta_\mu$ , on peut obtenir des estimations des effets individuels :

$$\hat{\alpha}_i = \frac{1}{T} \sum_{t=1}^T (Y_{it} - X'_{it}\hat{\beta}_\mu).$$

2. Régression quantile classique de la variable transformée  $\hat{Y}_{it} = Y_{it} - \hat{\alpha}_i$  sur les régresseurs  $X_{it}$ .

Canay montre également que ces estimateurs sont asymptotiquement normaux. Cependant, la matrice de variance-covariance ne correspond pas à celle produite par défaut dans la deuxième étape. Canay fournit un estimateur de la matrice de variance-covariance, mais propose également d'utiliser une procédure de *bootstrap* compte tenu de sa complexité. Dans cette procédure, il s'agit de répliquer, pour chacun des échantillons *bootstrap*, les deux étapes précédentes de l'estimation. Il n'existe pas encore dans les logiciels statistiques standard de procédure ou de package utilisant cette méthode mais on trouvera cependant sur le site de l'auteur un programme R permettant d'utiliser cette procédure en deux étapes<sup>13</sup>.

Canay n'établit la convergence de cet estimateur que lorsque le nombre de périodes  $T$  tend vers l'infini. Comme il le montre sur des données simulées, l'estimateur n'est pas toujours performant sur des panels contenant un petit nombre de périodes. D'autres méthodes ont également été proposées, toujours sous l'hypothèse (13) que les effets individuels ont un simple effet de translation. Koenker (2004) propose d'estimer simultanément les effets de  $X_{it}$  correspondant à  $q$

quantiles différents,  $(\beta_{\tau_k})_{k=1\dots q}$ , et les effets individuels  $(\alpha_i)_{i=1\dots n}$  (voir aussi Lamarche, 2010) :

$$\begin{aligned} & ((\hat{\beta}_{\tau_k})_{k=1\dots q}, (\hat{\alpha}_i)_{i=1\dots n}) \\ &= \arg \min_{\substack{(\alpha_i)_{i=1\dots n} \\ (\beta_{\tau_k})_{k=1\dots q}}} \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^q \rho_{\tau_k}(Y_{it} - X'_{it}\beta_{\tau_k} - \alpha_i) \\ &+ \lambda \sum_{i=1}^n |\alpha_i|. \end{aligned}$$

L'ajout du terme de pénalisation  $\lambda \sum_{i=1}^n |\alpha_i|$  permet d'éviter une trop grande dispersion des nombreux termes liés aux différents effets individuels. Ce coefficient de régularisation  $\lambda$  doit être fixé pour l'estimation : lorsqu'il est nul, aucune contrainte n'est donnée pour l'estimation des effets individuels et le risque existe d'obtenir des effets très dispersés. Lorsqu'il est très grand, la solution au programme de minimisation fournira des effets individuels estimés très faibles. Le package R *rqp* permet de faire cette estimation. Cette procédure est cependant très coûteuse en temps de calcul lorsque le nombre d'individus est élevé, et soulève la question du choix du coefficient de régularisation  $\lambda$ . Comme celui de Canay, cet estimateur n'est convergent qu'asymptotiquement avec  $T$ . Notons que l'hypothèse commune à ces deux estimateurs est que les effets individuels n'ont qu'un effet de translation sur la distribution d'intérêt, ce qui constitue une restriction importante. Des articles récents proposent l'utilisation d'estimateurs non-paramétriques, mais ces méthodes sont difficiles à mettre en œuvre et peu adaptées lorsque le nombre de variables explicatives est important. L'utilisation des données de panel pour des régressions quantiles est aujourd'hui encore un champ de recherche actif, et il est probable que d'autres estimateurs seront proposés dans les prochaines années.

### Les Quantile Treatment Effects

En pratique, on est souvent intéressé par l'effet non pas de l'ensemble des variables explicatives, mais plus spécifiquement de l'une d'entre elles. On peut par exemple s'intéresser à l'effet d'avoir suivi une formation professionnelle sur les revenus ou à l'effet d'une politique éducative sur la réussite scolaire... L'objectif est alors d'évaluer l'effet causal du fait d'avoir bénéficié du programme évalué, qu'on peut noter par une

13. L'adresse actuelle du site est <http://faculty.wcas.northwestern.edu/iac879/research.htm>.

indicatrice binaire  $T$ . Pour disposer d'un cadre pour traiter de ces questions, on peut utiliser le formalisme usuel de la littérature sur l'évaluation empirique des politiques publiques. Il consiste à dire que chaque personne a deux revenus potentiels,  $Y_0$  (celui qu'il peut espérer en l'absence du programme) et  $Y_1$  (celui qu'il peut espérer avec le programme). À ces revenus potentiels sont associées les deux distributions  $F_{Y_0}$  et  $F_{Y_1}$ . On peut alors définir le  $\tau^e$  *quantile treatment effect* (QTE) comme la distance horizontale entre les deux distributions (Doksum, 1974) :

$$\delta_\tau = q_\tau(Y_1) - q_\tau(Y_0).$$

De même, on peut définir sa restriction aux personnes qui ont effectivement bénéficié du programme (*quantile treatment effect on the treated*, QTET):

$$\delta_{\tau|T=1} = q_\tau(Y_1 | T=1) - q_\tau(Y_0 | T=1).$$

On peut ici s'attarder sur l'interprétation de ces paramètres. L'effet du traitement sur le quantile pour les traités (traduction très littérale du *quantile treatment effect on the treated*) correspond à la différence entre le  $\tau^e$  quantile de la distribution de la variable d'intérêt parmi les personnes qui ont bénéficié du programme  $T$ , et le quantile équivalent de la distribution de cette variable parmi ces mêmes personnes, si elles n'avaient pas bénéficié de ce programme. L'effet du traitement sur le quantile (*quantile treatment effect*) est plus général, puisqu'il correspond à la différence entre les quantiles des distributions qu'on s'attend à observer dans la population respectivement si le programme est généralisé à tous ou au contraire en son absence. Comme on l'a déjà souligné, ces paramètres ne correspondent pas, en général, à l'effet de ce programme pour les personnes qui se trouvent au niveau du  $\tau^e$  quantile de la distribution de  $Y$  en l'absence de ce programme. Ce paramètre ne permet donc pas en principe de dire si ce sont les personnes initialement les plus (dé)favorisées qui bénéficieraient du programme qu'on évalue. Pour passer à cette interprétation, il est à nouveau nécessaire de faire une hypothèse d'invariance des rangs, c'est-à-dire que les classement des individus selon les valeurs  $Y$  est le même en l'absence ou en présence du programme. Enfin, du fait de la non linéarité des quantiles, ces paramètres ne correspondent pas non plus à la distribution de l'effet du programme ( $q_\tau(Y_1) - q_\tau(Y_0) \neq q_\tau(Y_1 - Y_0)$ ). On trouvera dans Clements *et al.* (1997) une discussion des conditions sous lesquelles il est possible de borner cette distribution des effets.

Indépendamment de ces questions d'interprétation spécifiques aux régressions quantiles, la difficulté classique pour estimer ces effets du programme sur les quantiles tient à ce que, pour une personne donnée, on n'observe en fait que  $Y = TY_1 + (1-T)Y_0$ , c'est-à-dire le revenu potentiel avec traitement ( $Y_1$ ) si elle a bénéficié du programme et le revenu potentiel sans traitement sinon. On pourrait être tenté d'estimer simplement  $\delta_\tau$  (ou  $\delta_{\tau|T=1}$ ) par la différence de quantiles conditionnels des revenus observés,  $q_\tau(Y | T=1) - q_\tau(Y | T=0)$ . Mais, en général cette différence ne correspond pas au paramètre d'intérêt. En effet, dès qu'il existe une (auto) sélection dans l'entrée dans le programme (par exemple, lorsque les personnes qui ont choisi d'en bénéficier sont celles pour lesquelles il sera le plus efficace), la distribution des revenus observés parmi les bénéficiaires  $F_{Y|T=1}$  n'est pas représentative de la distribution du revenu potentiel avec application du programme à l'ensemble de la population  $F_{Y_1}$ .

Plusieurs méthodes ont été proposées pour identifier les effets moyens d'un programme en présence d'effets de sélection (on en trouvera par exemple une description dans Givord, 2010). Des extensions de ces méthodes à l'analyse des quantiles ont été proposées récemment. Nous ne présentons ici qu'une possibilité, qui est facilement implémentable. Elle repose sur l'hypothèse d'indépendance conditionnelle suivante (*Conditional Independence Assumption*, ou CIA) :

$$(Y_0, Y_1) \perp T | X. \quad (15)$$

Cette hypothèse correspond à l'exogénéité conditionnelle de  $T$  (on parle aussi de sélection sur observables, c'est-à-dire que toute la sélection dans le programme peut être expliquée par les variables observées  $X$ ). Elle est à la base des méthodes d'appariement (*matching*) ou simplement des régressions linéaires lorsque l'on s'intéresse à la moyenne. On pourrait donc envisager d'estimer l'impact du programme  $T$  par une régression quantile en contrôlant l'effet des observables  $X$ . Cependant, une telle régression quantile ne permet pas d'estimer directement le paramètre d'intérêt. Lorsque l'on inclut des variables de contrôles supplémentaires  $X$ , la régression quantile estime en effet le paramètre  $\hat{\delta}_\tau = q_\tau(Y_1 | X=x) - q_\tau(Y_0 | X=x)$ . Du fait de la non-linéarité des quantiles, ce paramètre ne correspond pas à  $\delta_\tau$  ni même à  $\delta_{\tau|T=1}$  en général. Les quantiles de la distribution des revenus potentiels  $Y_0$  et  $Y_1$  ne sont pas les mêmes que ceux des distributions de ces revenus potentiels conditionnelles aux observables.

Firpo (2007) propose une méthode pour résoudre ces deux problèmes. Celle-ci s'apparente aux méthodes d'appariement utilisées pour estimer l'effet moyen du traitement  $E(Y_1 - Y_0)$  sous l'hypothèse (15) d'indépendance conditionnelle. On fait tout d'abord une hypothèse de support commun, nécessaire également dans les méthodes d'appariement. Si  $p(X)$  est le score de propension à bénéficier du traitement, cette hypothèse s'écrit :

$$p(X) = P(T = 1 | X) \in ]0, 1[ \quad (16)$$

Cela signifie que pour chaque bénéficiaire du programme, il existe une personne qui n'en a pas bénéficié et présente les mêmes caractéristiques observables (et inversement). Firpo montre que sous les hypothèses ci-dessus il est possible d'identifier les deux quantiles  $q_\tau(Y_1)$  et  $q_\tau(Y_0)$ , à partir des seules données observées  $(Y, T, X)$ . Il utilise pour cela les relations :

$$\begin{aligned} \tau &= E \left[ \frac{T 1\{Y \leq q_\tau(Y_1)\}}{p(X)} \right] \\ &= E \left[ \frac{(1-T) 1\{Y \leq q_\tau(Y_0)\}}{1-p(X)} \right] \end{aligned} \quad (17)$$

qui découlent de (l'égalité est prouvée ici pour  $Y_1$ , le raisonnement étant identique pour  $Y_0$ ) :

$$\begin{aligned} E \left[ \frac{T 1\{Y \leq q_\tau(Y_1)\}}{p(X)} \right] &= E \left[ \frac{1\{Y_1 \leq q_\tau(Y_1)\}}{p(X)} E(T | Y_1, X) \right] \\ &= E \left[ \frac{1\{Y_1 \leq q_\tau(Y_1)\}}{p(X)} E(T | X) \right] \\ &= E[1\{Y_1 \leq q_\tau(Y_1)\}] \\ &= \tau \end{aligned}$$

Comme on observe  $(T, Y, X)$  et que l'on peut identifier le score de propension  $p(X)$ , ces relations permettent d'estimer  $q_\tau(Y_1)$  et  $q_\tau(Y_0)$ . En pratique, Firpo montre que l'on peut estimer  $\delta_\tau = q_\tau(Y_1) - q_\tau(Y_0)$  par une procédure en deux étapes :

1. estimer le score  $p(X)$ . Notons  $\hat{p}(X)$  un tel estimateur ;
2. estimer  $q_\tau(Y_0)$  puis  $q_\tau(Y_1)$  en utilisant une régression quantile sur la seule constante :

$$\hat{q}_\tau(Y_t) = \arg \min_b \sum \hat{\omega}_{t,i} \rho_\tau(Y_i - b) \quad (t = 0, 1),$$

$$\begin{aligned} \text{avec des pondérations } \hat{\omega}_{1,i} &= \frac{T_i}{\hat{p}(X_i)} \text{ et} \\ \hat{\omega}_{0,i} &= \frac{1 - T_i}{1 - \hat{p}(X_i)}. \end{aligned}$$

Il s'agit donc d'une régression quantile simple, mais pondérée afin de tenir compte des effets de sélection. Intuitivement, on pondère ainsi parmi les personnes non traitées celles qui ont néanmoins une probabilité plus grande de l'être. Firpo propose également un jeu de pondérations pour estimer l'effet du traitement sur les seuls traités  $\delta_{\tau|T=1}$  : dans ce cas on utilise comme pondération respectivement  $\hat{\omega}_{0,i|T=1} = \frac{\hat{p}(X_i)}{1 - \hat{p}(X_i)} \frac{(1 - T_i)}{\sum_{l=1}^n T_l}$  pour estimer  $q_\tau(Y_0 | T = 1)$  et  $\hat{\omega}_{1,i|T=1} = T_i / \sum_{l=1}^n T_l$  pour estimer  $q_\tau(Y_1 | T = 1)$ .

Cette méthode peut donc être implémentée simplement en utilisant des régressions quantiles standards, en pondérant les observations par le poids estimé correspondant au score de propension. Elle n'est en revanche valide que lorsque l'hypothèse d'indépendance conditionnelle aux observables est vérifiée, ce qui doit être discuté avec attention en fonction des modes de sélection dans le traitement. Lorsque cette hypothèse ne paraît pas plausible, il faut se tourner vers d'autres méthodes d'identification : selon le cas et les données disponibles, l'utilisation de variables instrumentales ou de données de panel sont des options possibles. Frandsen *et al.* (2012) proposent une extension aux régressions quantiles des méthodes de régression sur discontinuités (cas où l'affectation au traitement dépend d'une variable annexe de manière discontinue).

Au-delà du cadre strict de l'évaluation de politique publique, la méthode de Firpo (2007) est aussi utile lorsqu'on souhaite isoler l'effet d'une variable explicative binaire  $T$  sur la distribution d'une variable d'intérêt, en ayant contrôlé des autres caractéristiques observables  $X$ . Si l'on s'intéresse plutôt à une variable continue, on peut citer une méthode développée par ces auteurs (Firpo *et al.*, 2009), qui permet d'étudier l'effet d'une augmentation infinitésimale d'une seule variable explicative continue ( $X_2$  par exemple) sur la distribution de  $Y$ , les autres covariables ( $X_3, \dots, X_p$ ) restant inchangées. On parle dans ce cas de régression quantile « inconditionnelle » dans la mesure où l'objet d'intérêt est  $q_\tau(Y)$  (ou plus précisément sa variation suite à une variation infinitésimale de  $X_2$ ), plutôt que  $q_\tau(Y | X)$ .

## Les régressions quantiles dans les modèles non linéaires

Nous considérons enfin les extensions de la régression linéaire quantile aux modèles non-linéaires de la forme

$$Y = g(X'\beta_0 + \varepsilon), \quad (18)$$

où  $g$  est une fonction non linéaire connue. Deux exemples importants sont le modèle binaire, pour lequel  $g(x) = 1\{x > 0\}$ , et le modèle à censure fixe, pour lequel  $g(x) = \max(s, x)$  (ou  $g(x) = \min(s, x)$ ) avec  $s$  une constante connue. Ce dernier modèle est souvent utilisé pour modéliser la consommation d'un bien, qui prend la valeur nulle quand il n'est pas consommé (pour plus de détails, cf. par exemple Wooldridge, 2002, chap. 16). On suppose qu'il existe une variable latente  $c^*$  (liée à une optimisation d'utilité d'achat de ce bien par le consommateur éventuel), telle que la consommation observée vérifie  $c = \max(0, c^*)$ . Dans ces modèles, il est difficile d'utiliser des restrictions de la forme  $E(\varepsilon | X) = 0$  car en général,  $E(Y | X) \neq g(X'\beta_0)$ . L'approche standard consiste alors à imposer des hypothèses paramétriques sur la distribution des résidus. Par exemple, il est fréquent de supposer l'indépendance entre  $X$  et  $\varepsilon$  et la normalité de ces derniers (on parle alors de modèle probit lorsque  $g(x) = 1\{x > 0\}$  et de modèle tobit lorsque  $g(x) = \max(0, x)$ ). Ces hypothèses sont cependant restrictives et souvent difficiles à justifier.

Une approche alternative à ces hypothèses paramétriques est de recourir à des restrictions sur les quantiles. En effet, on peut facilement étendre les restrictions sur les quantiles des termes de perturbations  $\varepsilon$  à une transformation non linéaire, grâce à la propriété d'invariance déjà évoquée plus haut :

$$g(q_\tau(U)) = q_\tau(g(U)),$$

valable pour toute variable aléatoire  $U$  et toute fonction  $g$  croissante et continue à gauche<sup>14</sup>. Ainsi, si l'on impose dans le modèle non linéaire (18) la restriction  $q_\tau(\varepsilon | X) = 0$  et que  $g$  est croissante continue à gauche, on obtient :

$$q_\tau(Y | X) = g(q_\tau(X'\beta_0 + \varepsilon | X)) = g(X'\beta_0).$$

En reprenant un argument déjà utilisé plus haut, il s'ensuit que :

$$\beta_0 \in \arg \min_{\beta} E[\rho_\tau(Y - g(X'\beta))]$$

Comme précédemment, on estime alors  $\beta_0$  par :

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n \rho_\tau(Y_i - g(X_i'\beta)). \quad (19)$$

L'estimateur défini par (19) est très proche de celui de la régression quantile linéaire, la

différence étant simplement l'ajout dans le programme de la fonction  $g$ .

Pour le modèle de censure fixe pour lequel  $g(x) = \max(s, x)$ ,  $\hat{\beta}$  est  $\sqrt{n}$ -convergent, et peut être estimé par une application itérative de régressions quantiles linéaires (cf. par exemple Buchinsky, 1994). Pour la médiane ( $\tau = 1/2$ ), l'estimateur proposé par Powell (1984) (*censored LAD estimator*, i.e. l'estimateur des moindres valeurs absolues censuré) est implémenté sous Stata via la commande `clad`. Hong et Chernozhukov proposent par ailleurs un estimateur simple en trois étapes qui peut être utilisé pour l'ensemble des quantiles (cf. Fack et Landais, 2009, pour une utilisation de cette méthode sur données françaises).

## Exemples d'application

**P**our conclure cette présentation, et afin de familiariser le lecteur avec la présentation et la lecture des résultats d'une régression quantile, qui s'avèrent des exercices difficiles, nous proposons deux exemples de mise en œuvre de la méthode. Le premier est l'exemple très classique d'utilisation de la régression quantile pour l'analyse de la dispersion des salaires. Le second correspond à un usage plus avancé : l'évaluation des effets d'une politique publique, en présence d'un biais de sélection sur les bénéficiaires de cette politique.

### Comment lire les résultats d'une régression quantile ?

Pour le premier exemple, on ré-estime une équation de salaire classique à partir de l'enquête *Emploi 2008*. Cet exercice n'a d'autre prétention que d'illustrer les résultats issus d'une régression quantile sur un cas pratique, sans traiter le problème de la sélection dans l'emploi (voir sur ce point Buchinsky (1998)). Pour une étude plus complète de la question des rendements salariaux de l'expérience et de l'éducation, et de leur évolution en France, on se reportera par exemple à Charnoz *et al.* (2011).

La variable d'intérêt est le logarithme du salaire. Les variables explicatives sont les caractéris-

14. Une fonction  $g$  est continue à gauche si pour tout  $x$ , la limite de  $g(u)$  pour  $u$  tendant inférieurement vers  $x$  est égale à  $g(x)$ . Les fonctions  $g(x) = 1\{x > 0\}$  et  $g(x) = \max(0, x)$  sont donc continues à gauche.

tiques observables du salarié, à savoir le nombre d'années d'études, le sexe, sa nationalité, le nombre d'années d'expérience potentielle ainsi que le carré de celui-ci. Les estimations ont été faites pour chaque décile de la distribution conditionnelle du logarithme du salaire. On suppose donc, pour chaque décile :

$$\text{décile}_j(\ln(\text{salaire} | X)) = X' \beta_j$$

Les régressions quantiles permettent de déterminer comment varie chaque décile en fonction des déterminants auxquels on s'intéresse. Par exemple, le paramètre  $\beta_{k,j}$  dans la régression  $\text{décile}_j(\ln(\text{salaire}) | X) = X' \beta_j$  vérifie :

$$\beta_{k,j} = \text{décile}_j(\ln(\text{salaire}) | X_{-k}, X_k = x_k) - \text{décile}_j(\ln(\text{salaire}) | X_{-k}, X_k = x_k + 1).$$

Il s'agit du changement du  $j^{\text{e}}$  décile de la distribution conditionnelle de salaire suite à une augmentation d'une unité de  $X_k$ , par exemple une augmentation d'une année d'études, toutes choses égales par ailleurs (i.e., les autres variables  $X_{-k}$  restent constantes). Dans le cas d'une variable explicative binaire, comme le fait d'être un homme pour un salarié,  $\beta_{k,j}$  mesure simplement l'écart entre le  $j^{\text{e}}$  décile de la distribution des salaires des hommes (conditionnelle à l'ensemble des autres variables explicatives  $X_{-k}$ ) et le  $j^{\text{e}}$  décile de la distribution des salaires des femmes (également conditionnelle à  $X_{-k}$ ).

En termes de présentation, on remarquera qu'on a un jeu de coefficients estimés pour chaque quantile auquel on s'intéresse. Les résultats sont donc plus lourds à présenter que pour une régression classique. Dans la littérature, on les trouve présentés sous forme d'un tableau regroupant l'ensemble des coefficients, ou, de manière peut être plus parlante, sous forme de graphiques. C'est la solution que nous avons retenue ici (cf. graphique III).

Nous avons choisi de représenter les estimations des coefficients pour les différents déciles, avec l'intervalle de confiance à 95 % (zone en gris foncé), ainsi que, à titre de comparaison, la valeur et l'intervalle de confiance du coefficient des moindres carrés ordinaires (zone en gris clair).

Le coefficient correspondant à la constante peut être considéré comme le décile des salariés ayant les modalités de référence (ici, le fait d'être un

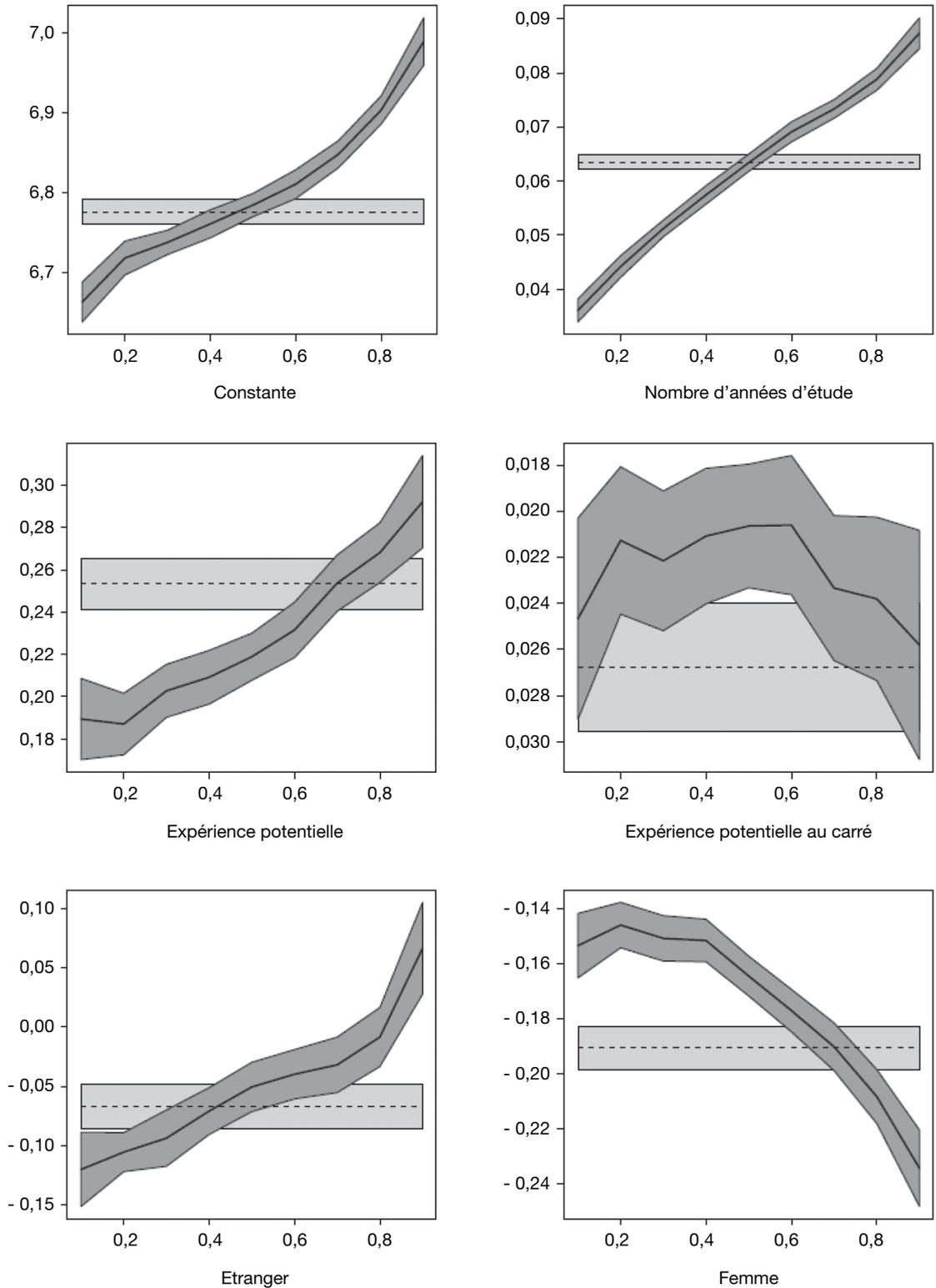
homme salarié avec la nationalité française, sans expérience potentielle et dont le niveau d'étude est minimal). Il est, sans surprise, croissant avec le décile (premier graphique en haut à gauche). On passe ainsi de 6,5 pour le premier décile à 7,0 pour le neuvième décile. Le coefficient estimé par les moindres carrés ordinaires, qui correspond au salaire moyen des salariés aux modalités de référence, sans expérience potentielle et de niveau d'étude minimal, est plus proche des premiers déciles (autour de 6,7), ce qui exprime bien le fait que la distribution du logarithme des salaires est asymétrique.

Le coefficient correspondant au nombre d'années d'étude est toujours positif, ce qui reflète le fait que le niveau d'étude décale globalement la distribution des salaires vers le haut. Son effet est très nettement croissant avec le décile. Le premier décile de la distribution conditionnelle des salaires augmente de 4 % quand le nombre d'années d'étude augmente d'une unité, une fois que l'on contrôle des autres variables observées, alors que la hausse est de 8 % pour le dernier décile. Une autre manière de présenter ce résultat est de dire que la dispersion des salaires augmente avec le nombre d'années d'études, ou encore que les distributions de salaires des plus diplômés sont plus inégales que celles des moins diplômés. C'est également le cas pour l'effet de l'expérience potentielle<sup>15</sup>.

Les salaires des femmes sont systématiquement inférieurs à ceux des hommes, mais ces différences sont d'autant plus fortes que l'on s'élève dans la distribution : conditionnellement aux autres caractéristiques observables, le neuvième décile de la distribution des salaires des femmes est ainsi inférieur de 24 % au neuvième décile de la distribution des salaires des hommes, tandis que cette différence n'est que de 15 % pour le premier décile. Ces différences peuvent s'expliquer par exemple par la présence de plafonds de verre qui compressent le haut de la distribution des salaires féminins. À l'inverse, le fait de ne pas disposer de la nationalité française a un impact négatif pour le bas de la distribution, mais les deux distributions conditionnelles se rapprochent ensuite. Le coefficient augmente avec le décile, il n'est plus significatif au niveau des septième et huitième déciles et il devient même positif au niveau du neuvième décile.

15. Du fait du terme quadratique, on peut interpréter l'augmentation marginale de l'expérience potentielle ( $\text{exppot}$ ) sur le décile  $j$  de la distribution de salaire comme  $b_{j1} + 2b_{j2} \text{exppot}$ , où  $b_{j1}$  est le coefficient relatif à  $\text{exppot}$  et  $b_{j2}$  le coefficient du carré d' $\text{exppot}$ .

Graphique III  
**Estimation des coefficients d'une équation de salaire par régressions quantiles.**



LECTURE : Le sixième décile du logarithme du salaire augmente de 0,07 lorsque le nombre d'années d'études augmente d'une unité, à valeurs données des autres variables explicatives. La hausse n'est que de 0,036 lorsqu'on s'intéresse au premier décile. La droite horizontale donne le coefficient de la même variable tel qu'obtenu par moindres carrés ordinaires. Il est de 0,064. Les bandes en grisé correspondent aux intervalles de confiance à 95 %.

Champ : population salariée, année 2008.

Source : enquête Emploi.

Ces résultats sont à interpréter avec la même prudence que ceux que donnerait la modélisation du salaire conditionnel moyen. La régression quantile est un outil qui permet d'estimer les effets de variables explicatives sur l'ensemble de la distribution d'une variable d'intérêt, mais elle ne règle aucun des éventuels problèmes d'endogénéité de certaines de ces variables. Par exemple, le fait d'avoir fait des études longues peut être lié à des compétences spécifiques ou à des réseaux familiaux qui ont également un effet positif sur le salaire. Ces éléments ne sont pas observés dans l'enquête. Dans ce cas, le coefficient des années d'études reflète également l'effet positif de ces caractéristiques inobservées, et pas uniquement l'effet causal d'une augmentation d'une année d'étude. De même, l'interprétation du coefficient correspondant au fait d'être étranger est délicate. Les coefficients négatifs puis positifs ne doivent pas forcément s'interpréter comme la présence d'une discrimination négative puis positive envers les étrangers. Il faudrait supposer que la distribution des caractéristiques ayant un effet sur le salaire de l'ensemble des salariés étrangers travaillant en France est identique à cette même distribution pour les salariés français. Or les étrangers décidant de travailler en France ont sans doute des profils professionnels particuliers, qui peuvent expliquer ces coefficients. Au final, exactement les mêmes précautions d'interprétation s'imposent que dans le cas d'une régression linéaire classique. En cas d'endogénéité de certaines variables explicatives, il est nécessaire, pour obtenir une interprétation causale des coefficients, de mobiliser par exemple la méthode des variables instrumentales décrite précédemment, toute la difficulté étant de trouver des instruments adéquats.

D'autre part, les estimations obtenues correspondent à l'effet des variables explicatives sur les distributions de la variable d'intérêt conditionnelles à ces variables. Elles renseignent sur les écarts entre les quantiles d'ordre  $\tau$  de la distribution des salaires conditionnelle à ces variables. Comme expliqué plus haut en encadré, elle ne permet pas d'évaluer directement comment se modifierait le quantile d'ordre  $\tau$  de la distribution de salaires de l'ensemble de la population si la distribution de ces variables explicatives était différente, par exemple si la proportion de diplômés du supérieur était supérieure. Ceci vient de la propriété de non-linéarité des quantiles déjà évoquée : les quantiles de la population entière ne s'obtiennent pas simplement en intégrant les quantiles conditionnels qui sont modélisés dans les régressions quantiles.

Enfin, comme expliqué également plus haut, les résultats des régressions quantiles, tout comme ceux obtenus par une régression linéaire, n'ont pas directement d'interprétation individuelle. Le principe des régressions quantiles est *stricto sensu* de comparer des distributions conditionnelles entre elles. Elles permettent par exemple de dire que le premier décile des salaires des salariés étrangers est inférieur de 12 % à celui des salariés ayant la nationalité française, toutes choses égales par ailleurs. En revanche, elles ne permettent pas a priori de dire que le salaire d'un salarié étranger qui se trouve au niveau du premier décile de la distribution de salaire de cette population augmenterait d'autant s'il acquérait la nationalité française. Pour pouvoir ainsi interpréter les résultats obtenus, il faut supposer que ce salarié occuperait le même rang dans les deux distributions (correspondant respectivement aux salaires des salariés de nationalité française et aux salaires des salariés ne disposant pas de la nationalité française). Cette hypothèse d'invariance des rangs a le mérite de fournir une interprétation simple des coefficients des régressions quantiles mais elle est restrictive. Cette mise en garde dans l'interprétation des résultats est de même nature que celle qu'on peut avoir pour les résultats de la modélisation de la moyenne du logarithme des salaires par une régression linéaire. Les estimations ainsi obtenues montrent que le salaire moyen des étrangers est inférieure de 7 % à celui des salariés français. Cela ne signifie pas qu'un salarié étranger dont le salaire se trouve au niveau du salaire moyen verrait son salaire augmenter d'autant s'il acquérait la nationalité française.

### Un exemple de régression quantile instrumentale

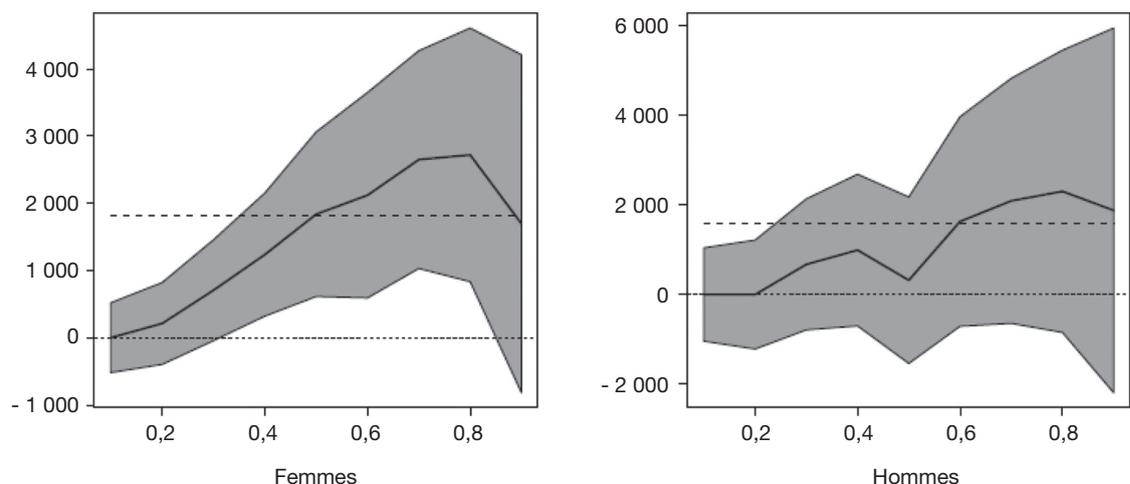
Notre deuxième exemple reprend l'application d'Abadie *et al.* (2002). Il permet d'illustrer l'utilisation des méthodes de régressions quantiles instrumentales présentée en troisième partie. Il utilise les données issues de l'expérimentation de l'efficacité d'un programme de formation de chômeur, le *Job Training Partnership Act (JTPA)*, mis en place à partir de 1983 aux États-Unis. Il s'agit d'un ensemble de programmes de formation et d'assistance destinés aux jeunes défavorisés. L'évaluation de l'efficacité de ce genre de programme est souvent rendue difficile par les effets d'auto-sélection : en général, ce sont les personnes qui peuvent en retirer le plus grand bénéfice qui choisissent de rentrer dans le dispositif.

Pour contrôler ce biais d'auto-sélection, une expérimentation a été mise en place entre 1987 et 1989 dans 16 structures locales auprès d'un échantillon initial de 20 000 jeunes environ. Les programmes de formation correspondant au JTPA n'ont été proposés qu'à deux tiers de ces jeunes, tirés aléatoirement. Des données ont ensuite été collectées sur la séquence de revenus de l'ensemble des jeunes de l'échantillon initial. Mais cette procédure de sélection aléatoire n'est pas suffisante pour qu'on puisse directement comparer les jeunes ayant participé au programme de formation à ceux qui n'y ont pas participé. Cela tient au fait que les personnes affectées au programme par tirage au sort pouvaient néanmoins choisir de ne pas en profiter. Ainsi, seules 60 % d'entre elles ont profité des programmes de formation. Inversement, l'entrée dans le programme n'était pas complètement fermée aux jeunes qui n'avaient pas été sélectionnés par le tirage au sort puisque 2 % d'entre eux les ont néanmoins suivis. Ainsi, malgré le tirage au sort, il y a une part d'auto-sélection dans la participation effective au programme. Mais on dispose en revanche d'un instrument, l'affectation aléatoire au dispositif. Issue d'un tirage au sort, elle n'est évidemment pas corrélée aux déterminants inobservés du revenu. En revanche, elle explique fortement le fait d'avoir bénéficié ou non du programme.

Nous appliquons donc la méthode de Chernozhukov et Hansen (2008) décrite précédemment pour évaluer l'impact de ce programme sur l'ensemble des jeunes<sup>16</sup>. Les estimations sont conduites séparément pour les hommes et pour les femmes. On a estimé ici l'ensemble des déciles. Les résultats sont présentés sur le graphique IV. Les régressions instrumentales fournissent des résultats près de deux fois plus faibles que les simples régressions linéaires (droites pointillées en gras), ce qui traduit bien l'auto-sélection dans le dispositif. Les régressions quantiles indiquent également que la moyenne masque de grandes disparités dans l'effet du programme. Pour les femmes, l'effet moyen du traitement, c'est-à-dire la différence entre la moyenne espérée des revenus en l'absence du programme et celle des revenus avec le programme est de 1 825 dollars. Mais il n'augmente en fait que de 210 dollars le deuxième décile de la distribution de revenus, alors que l'augmentation atteint 2 720 dollars pour le huitième décile. Les différences sont également très marquées pour les hommes, mais les estimations sont bien plus imprécises et ne permettent jamais d'exclure leur nullité aux seuils ordinaires de significativité. Les intervalles de confiance sont en effet larges. De manière similaire aux régressions linéaires classiques, la précision

16. Les données sont disponibles à l'adresse <http://econ-www.mit.edu/faculty/angrist/data1/data/abangim02>.

Graphique IV  
Estimation de l'impact d'un programme de formation, régression quantile instrumentée.



Lecture : les paramètres estimés sont la moyenne (en pointillé, estimation par les doubles moindres carrés) et les neuf déciles. Les zones grisées correspondent à l'intervalle de confiance à 95 %, estimé par une méthode de bootstrap. L'effet moyen de la formation sur les revenus des jeunes femmes est de 1 825 dollars. L'impact de cette formation sur le premier décile de la distribution des revenus des jeunes femmes est proche de zéro.

Champ : bénéficiaires et groupe témoin du Job Training Partnership Act, États-Unis, 1983.

Source : calcul des auteurs à partir des données de Abadie et al., (2002).

des paramètres estimés à partir d'une méthode instrumentale est bien plus faible qu'avec une régression quantile standard, même quand les instruments prédisent correctement les variables explicatives endogènes, comme c'est le cas ici.

En termes d'interprétation, on rappelle que ces valeurs correspondent à la différence des quantiles des distributions de revenus que l'on s'attend à observer respectivement en l'absence ou en présence du programme de formation. Cela ne signifie pas a priori que les femmes dont le revenu serait au niveau du deuxième décile en l'absence de la formation vont bénéficier d'une augmentation de 210 dollars grâce à celle-ci, sauf à faire l'hypothèse que la formation ne modifie pas l'ordre relatif des revenus des personnes. Or ceci peut ne pas être le cas si le programme est plus profitable aux personnes situées en bas de l'échelle des revenus qu'aux personnes situées au dessus, auquel cas les niveaux des revenus potentiels peuvent s'invertir. Par ailleurs, alors qu'on peut interpréter le résultat obtenu par les doubles moindres carrés comme l'effet moyen de la formation (et pas seulement comme la différence des moyennes des revenus avec et sans la formation), il n'est pas en général possible d'interpréter des résultats obtenus par la régression quantile en termes de distribution des effets de la politique. Autrement dit, il n'est pas possible de dire qu'un cinquième des jeunes femmes va bénéficier d'au moins 210 dollars grâce à la formation, car la différence des premiers quartiles des distributions de revenus avec et sans formation ne correspond pas a priori au premier quartile des gains individuels liés à la formation. Les régressions quantiles permettent de décrire comment se déforme la distribution des niveaux de revenus (quelle nouvelle répartition des richesses on observe grâce au programme), mais pas la distribution des évolutions individuelles (comment se répartissent les gains). Cela n'en réduit pas l'intérêt. Si le gain moyen est une information essentielle pour l'évaluation, en regard par exemple des coûts engagés pour ce programme, la répartition finale des niveaux de revenus est également importante, dans la mesure où elle permet notamment de déterminer si au moins  $x$  % des personnes sont au-dessus d'un seuil minimal de revenu après application de la politique. Cette information est souvent plus pertinente que l'identification précise des gagnants et des perdants.

\* \*  
\*

En conclusion, la régression quantile est un outil facile d'utilisation, qui permet d'enrichir la description quantitative des phénomènes économiques et sociaux. Cet article en a présenté les principes et a cherché à préciser l'utilisation qui peut en être faite. On insiste sur le fait que les problèmes classiques d'endogénéité éventuelle des variables explicatives doivent être traités, ce qui peut nécessiter de mobiliser des outils d'identification spécifiques détaillés dans la troisième section (à condition de disposer des instruments ou des données nécessaires). Par ailleurs, les régressions quantiles modélisent les quantiles conditionnels de la distribution d'intérêt. L'objet d'étude est donc avant tout cette distribution. Tout comme l'objet d'une régression linéaire n'est pas de renseigner sur l'effet de covariables sur un hypothétique « individu moyen », le premier objectif des régressions quantiles n'est pas de décrire l'effet de ces covariables sur des individus « assignés » à une place dans la distribution. Par exemple, on pourra déduire de l'évaluation des effets d'un programme éducatif par une régression quantile qu'il permet d'augmenter de  $x$  points le niveau atteint par au moins 90 % des élèves. En toute rigueur, on ne peut dire sans hypothèse supplémentaire qu'il permet d'augmenter d'autant le niveau des élèves les plus faibles en l'absence de ce programme. Cette interprétation suppose que les enfants les plus en difficulté avec une certaine méthode d'apprentissage le sont également face à une autre méthode. Cette hypothèse d'invariance des rangs permet de tirer davantage d'information des résultats de la régression. Mais elle n'est pas indispensable pour mettre en œuvre la méthode et en déduire des résultats utiles.

Expliciter les précautions indispensables à l'interprétation de la régression quantile ne doit donc pas conduire le lecteur à conclure qu'elle est un outil sophistiqué mais peu utile pour l'analyse. Le bilan d'un programme ne se juge pas sur son effet pour tel ou tel individu, mais sur sa capacité à améliorer ou non une situation générale. Ainsi, estimer qu'un programme augmente significativement les plus bas déciles mais n'a pas d'effet sur les déciles supérieurs signifie qu'il a permis de réduire les inégalités de résultats, sans diminution des ambitions scolaires (les résultats des têtes de classe sont aussi bons, même si, en principe, rien ne garantit que les têtes de classes sont les mêmes avec et sans programme). Quantiles et fonction de répartition étant naturellement liés, on peut aussi, au prix de quelques manipulations, tirer des conclusions du fait que le programme permet de ramener de tant à tant la proportion d'élèves

en dessous d'un seuil (dans notre exemple, un niveau minimum d'acquisition des connaissances). De fait, de nombreux articles récents ont montré que l'analyse au-delà de la moyenne permet d'enrichir les évaluations quantitatives de politiques publiques. Ainsi Bitler *et al.* (2006) montrent qu'une même réforme peut avoir des effets négligeables en moyenne tout en modifiant significativement la distribution des revenus. Casalone et Sonedda (2013) évaluent les effets d'une réforme fiscale sur la distribution des revenus, Jackson et Page (2013) l'impact d'un programme de réduction de la taille des classes (projet STAR) sur la réussite

scolaire des enfants (voir aussi Bitler *et al.*, 2014). Indépendamment de leur mobilisation pour l'évaluation des politiques publiques, les régressions quantiles sont également des outils adaptés à l'étude de l'évolution des inégalités de revenus, dans le sillage des travaux de Buchinsky sur le sujet. Enfin, même lorsque l'objet d'intérêt n'est pas l'ensemble de la distribution, les propriétés statistiques des régressions quantiles en font une alternative utile aux méthodes économétriques classiques, qu'il s'agisse de traitement de variables tronquées ou censurées ou de garantir une meilleure robustesse en présence de valeurs aberrantes. □

---

## BIBLIOGRAPHIE

- Abadie A., Angrist J. et Imbens G. (2002)**, « Instrumental variables estimates of the effect of subsidized training on the quantiles of trainee earnings », *Econometrica*, vol. 70, n° 1, pp. 91-117.
- Aeberhardt R., Fougère D., Pouget J. et Rathelot R. (2010)**, « L'emploi et les salaires des enfants d'immigrés », *Économie et Statistique*, n° 433, pp. 31-46.
- Biliyas Y. et Koenker R. (2001)**, « Quantile regression for duration data : A reappraisal of the Pennsylvania reemployment bonus experiments », *Empirical Economics*, vol. 26, n° 1, pp. 199-220.
- Bitler M.P. Gelbach J.B. et Hoynes H.W. (2006)**, « What Mean Impacts Miss: Distributional Effects of Welfare Reform Experiments », *American Economic Review*, vol. 96, n° 4, pp. 988-1012.
- Bitler M. P., Gelbach J. B. et Hoynes H. W. (2014)**, « Can Variation in Subgroups Average Treatment Effects Explain Treatment Effect Heterogeneity ? Evidence from a Social Experiment », *NBER Working Papers* n° 20142, National Bureau of Economic Research, Inc..
- Biscourp P., Boutin X. et Vergé T. (2013)**, « The Effects of Retail Regulations on Prices: Evidence from the Loi Galland », *The Economic Journal*, Vol. 123 (12), pp. 1279-1312.
- Buchinsky M. (1994)**, « Changes in the U.S. Wage Structure 1963-1987: Application of Quantile Regression », *Econometrica*, vol. 62, n° 2, pp. 405-58.
- Buchinsky M. (1998)**, « The dynamics of changes in the female wage distribution in the USA : a quantile regression approach », *Journal of Applied Econometrics*, vol. 13, n° 1, pp. 1-30.
- Cade B.S. et Noon B.R. (2003)**, « A Gentle Introduction to Quantile Regression for Ecologists », *Frontiers in Ecology and The Environment*, vol. 1, pp. 412-420.
- Canay I.A. (2011)**, « A simple approach to quantile regression for panel data », *The Econometrics Journal*, vol. 14, n° 3, pp. 368-386.
- Casalone G. et Sonedda D. (2013)**, « Evaluating The Distributional Effects Of Fiscal Policies Using Quantile Regressions », *Review of Income and Wealth*, vol. 59, n° 2, pp. 305-325.
- Charnoz P., Coudin E. et Gaini M. (2011)**, « Wage inequalities in France 1976-2004: a quantile regression analysis », *Document de Travail Insee- DESE*, n° g2011-06.
- Chernozhukov V., Fernandez-Val I. et Galichon A. (2010)**, « Quantile and Probability Curves Without Crossing », *Econometrica*, vol. 78, n° 3, pp. 1093-1125.
- Chernozhukov V. et Hansen C. (2008)**, « Instrumental variable quantile regression: A robust inference approach », *Journal of Econometrics*, n° 142, pp. 379-398.
- Clements N., Heckman J. et Smith J. (1997)**, « Making the most out of programme evaluations and social experiments: Accounting for heterogeneity in programme impacts », *Review of Economic Studies*, vol. 64, n° 4, pp. 487-535.

- Corneec M. (2014)**, « Constructing a conditional GDP fan chart with an application to french business survey data », *Journal of Business Cycle Measurement and Analysis*, vol. 2013/2, n° 9, pp. 109-127.
- Doksum K. (1974)**, « Empirical probability plots and statistical inference for nonlinear models in the two-sample case », *The Annals of Statistics*, vol. 2, n° 2, pp. 267-277.
- Dumontet M. et Franc C. (2014)**, « Gender differences in French GPs' activity: the contribution of quantile regressions », *European Journal of Health Economics*. À paraître.
- Etienne J.M. et Narcy M. (2010)**, « Gender Wage Differentials in the French Nonprofit and For-Profit Sectors: Evidence from Quantile Regression », *Annales d'Économie et Statistique*, n° 99/100, pp. 67-90.
- Fack G. et Landais C. (2009)**, « Les incitations fiscales aux dons sont-elles efficaces ? », *Économie et Statistique*, n° 427, pp. 101-121.
- Firpo S. (2007)**, « Efficient semiparametric estimation of quantile treatment effects », *Econometrica*, vol. 75, n° 1, pp. 259-276.
- Firpo S., Fortin N. M. et Lemieux T. (2009)**, « Unconditional quantile regressions », *Econometrica*, vol. 77, n° 3, pp. 953-973.
- Fitzenberger B. et Wilke R. A. (2005)**, « Using quantile regression for duration analysis », *ZEW Discussion Papers 05-65*, ZEW / Center for European Economic Research.
- Fortin N. Lemieux T. et Firpo S. (2011)**, « Decomposition Methods in Economics », volume 4 du *Handbook of Labor Economics*, chapitre 1, pp. 1-102. Elsevier.
- Frandsen B.R., Frolich M. et Melly B. (2012)**, « Quantile treatment effects in the regression discontinuity design », *Journal of Econometrics*, vol. 168, n° 2, pp. 382-395.
- Givord P. (2010)**, « Méthodes économétriques pour l'évaluation des politiques publiques », *Document de Travail Insee- DESE*, n° G2010-08.
- Givord P. et D'Haultfoeuille X. (2013)**, « La régression quantile en pratique », *Document de travail Insee-DMCSI*, n° L2013-1
- He X. et Hu F. (2002)**, « Markov chain marginal bootstrap », *Journal of the American Statistical Association*, vol. 97, n° 459, pp. 783-795.
- Hong H. et Chernozhukov V. (2002)**, « Three-step censored quantile regression and extramarital affairs », *Journal of the American Statistical Association*, vol. 97, pp. 872-882.
- Hyndman R.J. et Fan Y. (1996)**, « Sample quantiles in statistical packages », *The American Statistician*, vol. 50, n° 4, pp. 361-365.
- Jackson E. et Page M.E. (2013)**, « Estimating the distributional effects of education reforms: A look at Project STAR », *Economics of Education Review*, vol. 2 (C), pp. 92-103.
- Kocherginsky M., He X. et Mu Y. (2005)**, « Practical confidence intervals for regression quantiles », *Journal of Computational and Graphical Statistics*, vol. 14, n° 1, pp. 41-55.
- Koenker R. et Machado J.A.F. (1999)**, « Goodness of fit and related inferences processes for quantile regression », *Journal of the American Statistical Association*, vol. 94, n° 448, pp. 1296-1310.
- Koenker R. et Hallock K.F. (2001)**, « Quantile regression », *Journal of Economic Perspectives*, vol. 15, n° 4, pp. 143-156.
- Koenker R. (2004)**, « Quantile regression for longitudinal data », *Journal of Multivariate Analysis*, vol. 91, n° 1, pp. 74-89.
- Koenker R. (2005)**, *Quantile Regression*. Econometric Society Monograph Series, Cambridge University Press.
- Lamarche C. (2010)**, « Robust penalized quantile regression estimation for panel data », *Journal of Econometrics*, vol. 157, n° 2, pp. 396-408.
- Machado J. A. et Silva J. M. C. S. (2005)**, « Quantiles for counts », *Journal of the American Statistical Association*, vol. 100, pp. 1226-1237.
- Meurs D. et Ponthieux S. (2006)**, « L'écart des salaires entre les femmes et les hommes peut-il encore baisser ? », *Économie et Statistique*, n° 398, pp. 99-129.
- Portnoy S. (1992)**, « Asymptotic behavior of the number of regression quantile breakpoints », *SIAM Journal on Scientific and Statistical Computing*, vol. 12, n° 4, pp. 867-883.
- Portnoy S. et Koenker R. (1997)**, « The gaussian hare and the laplacian tortoise : Computability of squared- error versus absolute-error estimators », *Statistical Science*, vol. 12, n° 4, pp. 279-296.

**Powell J.L. (1991)**, *Estimation of monotonic regression models under quantile restrictions*. Cambridge: Cambridge University Press.

**Powell J. (1984)**, « Least absolute deviations estimation for the censored regression model », *Journal of Econometrics*, vol. 25, pp. 303-325.

**Samson A.-L. (2006)**, « La dispersion des honoraires des omnipraticiens sur la période 1983-2004. Une application de la méthode des régressions quantiles », *Document de travail de la Drees, série Etudes*, n° 62.

**Wooldridge J. W. (2002)**, *Econometric Analysis of Cross Section and Panel Data*, MIT Press.

---

PROPRIÉTÉS DES QUANTILES ET DÉTAILS SUR L'INFÉRENCE

1 - Définitions

Le quantile d'ordre  $\tau \in (0,1)$  d'une variable aléatoire réelle  $U$  est défini par :

$$q_\tau(U) = \inf\{x \mid F_U(x) \geq \tau\},$$

$F_U$  étant la fonction de répartition de  $U$ . Dans le cas où  $F_U$  est continue et strictement croissante, on a simplement  $q_\tau(U) = F_U^{-1}(\tau)$ . Le graphique ci-dessous illustre la définition des quantiles dans le cas général.

Pour deux variables aléatoires  $U$  et  $V$ , le quantile conditionnel  $q_\tau(U \mid V)$  est défini de manière similaire par :

$$q_\tau(U \mid V) = \inf\{x \mid F_{U|V}(x) \geq \tau\},$$

où  $F_{U|V}$  est la fonction de répartition de  $U$  conditionnelle à  $V$ .

2 - Robustesse aux valeurs aberrantes

Nous démontrons tout d'abord les propriétés de robustesse des régressions quantiles aux valeurs aberrantes énoncées dans le texte. On suppose tout d'abord

$$Y = AX'\delta + (1 - A)(X'\gamma + \varepsilon),$$

avec  $\varepsilon, A$  et  $X$  mutuellement indépendants.  $A$  est une variable inobservée valant 1 lorsque  $Y$  est aberrant, 0 sinon et  $P(A = 1) = p$ . Le paramètre  $\beta_\tau$  de la régression quantile vérifie, pour tout  $x$ ,  $\tau = P(Y \leq x'\beta_\tau \mid X = x)$ . Par conséquent,

$$\begin{aligned} \tau &= P(Y \leq X'\beta_\tau \mid X = x) \\ &= P(A = 1 \mid X = x)P(Y < X'\beta_\tau \mid X = x, A = 1) \\ &\quad + P(A = 0 \mid X = x)P(Y \leq X'\beta_\tau \mid X = x, A = 0) \\ &= pP(x'\delta < x'\beta_\tau \mid X = x) \\ &\quad + (1 - p)P(x'\gamma + \varepsilon \leq x'\beta_\tau \mid X = x, A = 0) \\ &= (1 - p)P(\varepsilon < x'(\beta_\tau - \gamma)) \end{aligned}$$

où la dernière égalité vient du fait que  $x'\delta$  est supposé très grand, et donc supérieur à  $x'\beta_\tau$ , et  $\varepsilon$  est indépendant de  $(A, X)$ . Ainsi, pour tout  $x$ ,

$$x'(\beta_\tau - \gamma) = q_{\tau/(1-p)}(\varepsilon)$$

On en déduit que  $\beta_{k,\tau} = \gamma_k$  pour tout  $k > 1$  et  $\beta_{1,\tau} = \gamma_1 + q_{\tau/(1-p)}(\varepsilon)$ . Si l'on excepte la constante, les coefficients de la régression quantile que l'on obtient sont donc égaux à ceux que l'on obtiendrait en l'absence de valeur aberrante.

Considérons maintenant le modèle plus général suivant :

$$Y = AX'\delta + (1 - A)(X'\beta_U),$$

où  $U, A$  et  $X$  sont mutuellement indépendants et  $U$  suit une loi uniforme sur  $[0,1]$ . Notons  $\hat{\beta}_\tau$  le coefficient de la régression quantile de  $Y$  sur  $X$ . En suivant le même raisonnement que précédemment, on obtient l'équation suivante en  $\hat{\beta}_\tau$ , valable pour tout  $x$  :

$$\tau = (1 - p)P(x'\beta_U < x'\hat{\beta}_\tau). \tag{20}$$

Le paramètre  $\beta_{\tau/(1-p)}$  est solution de cette équation. Dès que le modèle est identifié, il existe une seule solution vérifiant (20) pour tout  $x$ . Par conséquent,  $\hat{\beta}_\tau = \beta_{\tau/(1-p)}$ .

3 - Invariance à une transformation monotone

Les quantiles satisfont l'importante propriété d'invariance suivante.

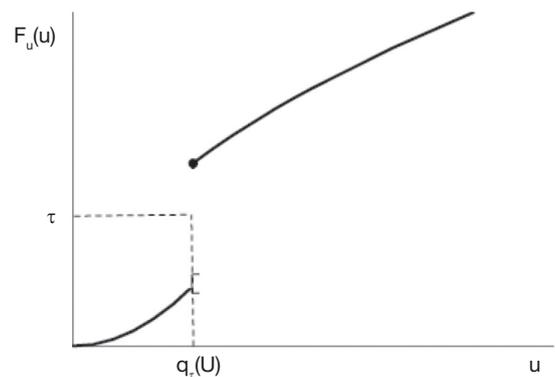
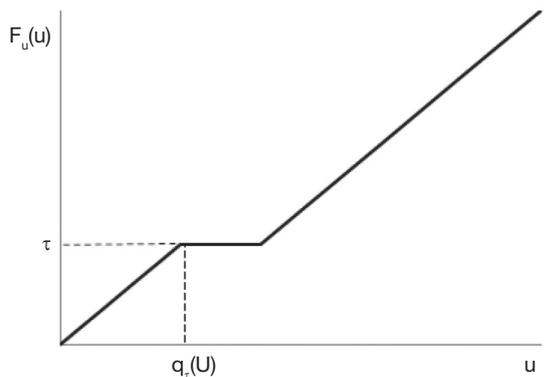
**Proposition 1 :** Soit  $g$  une fonction croissante et continue à gauche. Alors :

$$g(q_\tau(U)) = q_\tau(g(U)).$$

**Démonstration :** Grâce à la monotonie de  $g$  on a  $P(U \leq q_\tau(U)) = P(g(U) \leq g(q_\tau(U)))$  et par définition de  $q_\tau(U)$  :  $\tau \leq P(U \leq q_\tau(U))$ . Or par définition on a aussi  $q_\tau(g(U)) = \inf\{x \in \mathbb{R} \mid F_{g(U)}(x) \geq \tau\}$ , donc  $g(q_\tau(U)) \geq q_\tau(g(U))$ . Réciproquement, en définissant  $g^-(v) = \sup\{x \mid g(x) \leq v\}$ , on a :

$$P(g(U) \leq q_\tau(g(U))) \leq P(U \leq g^-(q_\tau(g(U))))$$

Graphique  
Quantile d'une variable dans le cas général



Par définition de  $q_\tau(g(U))$  et  $q_\tau(U) = \inf\{x \in \mathbb{R} \mid F_U(x) \geq \tau\}$  on en déduit que:  $g^-(q_\tau(g(U))) \geq q_\tau(U)$ . De la continuité de  $g$  à gauche, on a aussi que  $g(g^-(q_\tau(g(U)))) \leq q_\tau(g(U))$ . Donc  $q_\tau(g(U)) \geq g(q_\tau(U))$ , ce qui achève la démonstration.

Ce résultat implique notamment que  $q_\tau(aU + b) = aq_\tau(U) + b$ , ou, de même,  $q_\tau(a(X)U + b(X) \mid X) = a(X)U + b(X)$ . Mais il implique également que  $q_\tau(\max(s, U)) = \max(s, q_\tau(U))$ , ou que  $q_\tau(\mathbf{1}\{U > 0\}) = \mathbf{1}\{q_\tau(U) > 0\}$ . En revanche, et contrairement à l'espérance, la fonction quantile n'est pas linéaire : on a en général  $q_\tau(U_1 + U_2) \neq q_\tau(U_1) + q_\tau(U_2)$ .

#### 4 - Quantiles et minimisation de perte

La propriété suivante est cruciale pour l'estimation.

**Proposition 2 :** Soit  $\rho_\tau(u) = (\tau - \mathbf{1}\{u < 0\})u$ . On a :

$$q_\tau(U) \in \underset{a}{\operatorname{argmin}} E[\rho_\tau(U - a)].$$

**Démonstration :** Soit  $m(a) = E[\rho_\tau(U - a)]$ . Prenons  $a \leq q_\tau(U)$  et montrons que  $m(a) \geq m(q_\tau(U))$ . On a, après quelques calculs,

$$m(a) - m(q_\tau(U)) = (q_\tau(U) - a)\tau + aF_U(a) - q_\tau(U)F_U(q_\tau(U)^-) + E[U\mathbf{1}\{U \in [a, q_\tau(U)]\}],$$

où  $F_U(q_\tau(U)^-) = \lim_{x \uparrow q_\tau(U)} F_U(x)$ . Par ailleurs,

$$E[U\mathbf{1}\{U \in [a, q_\tau(U)]\}] \geq aE[\mathbf{1}\{U \in [a, q_\tau(U)]\}] = a(F_U(q_\tau(U)^-) - F_U(a)).$$

Donc

$$m(a) - m(q_\tau(U)) \geq (q_\tau(U) - a)(\tau - F_U(q_\tau(U)^-)).$$

Par définition des quantiles,  $F_U(q_\tau(U)^-) \leq \tau$ . Donc  $m(a) \geq m(q_\tau(U))$ . On montre de même que  $m(a) \leq m(q_\tau(U))$  pour tout  $a \geq q_\tau(U)$ .

Notons que le minimum de  $a \mapsto E[\rho_\tau(U - a)]$  n'est pas unique en général. Ceci provient du fait que l'équation  $F_U(a) = \tau$  peut avoir plusieurs solutions. Le minimum est cependant unique si la fonction de répartition de  $U$  est strictement croissante.

#### 5 - Inférence par estimation directe

Cette approche consiste à estimer directement la variance asymptotique en partant de la formule (9) du texte. Dans le cas général, la difficulté principale est d'estimer  $J_\tau = E(f_{\varepsilon_\tau} \mid X(0) \mid XX')$ . Pour ce faire, Powell (1991) propose de s'appuyer sur l'idée suivante :

$$J_\tau = \lim_{h \rightarrow 0} E \left[ \frac{\mathbf{1}\{|\varepsilon_\tau| \leq h\}}{2h} XX' \right].$$

On estime alors  $J_\tau$  par

$$\hat{J}_\tau = \frac{1}{2nh_n} \sum_{i=1}^n \mathbf{1}\{|\hat{\varepsilon}_{i\tau}| \leq h_n\} X_i X_i'. \quad (21)$$

où  $h_n \rightarrow 0$  et  $\sqrt{nh_n} \rightarrow \infty$ .

Cette formule est plus simple dans le cas du modèle de translation, puisque seule l'estimation de  $1/f_\varepsilon(q_\tau(\varepsilon))$  est problématique. Soit  $\hat{\varepsilon}_{i\tau} = Y_i - X_i' \hat{\beta}_\tau$ . On a :

$$\begin{aligned} \frac{1}{f_\varepsilon(q_\tau(\varepsilon))} &= \frac{1}{f_\varepsilon(F_\varepsilon^{-1}(\tau))} \\ &= \frac{\partial F_\varepsilon^{-1}}{\partial \tau}(\tau) \\ &= \lim_{h \rightarrow 0} \frac{F_\varepsilon^{-1}(\tau + h) - F_\varepsilon^{-1}(\tau - h)}{2h} \end{aligned}$$

On peut donc estimer  $1/f_\varepsilon(q_\tau(\varepsilon))$  par

$$\frac{\hat{F}_\varepsilon^{-1}(\tau + h_n) - \hat{F}_\varepsilon^{-1}(\tau - h_n)}{2h_n},$$

où  $\hat{F}_\varepsilon^{-1}$  est le quantile empirique de  $\hat{\varepsilon}_T$ . L'estimateur de la variance asymptotique vaut alors :

$$\hat{V}_{as} = \tau(1-\tau) \left( \frac{\hat{F}_\varepsilon^{-1}(\tau + h_n) - \hat{F}_\varepsilon^{-1}(\tau - h_n)}{2h_n} \right)^2 \left[ \frac{1}{n} \sum_{i=1}^n X_i X_i' \right]^{-1}. \quad (22)$$

Cet estimateur est parfois proposé par défaut dans des logiciels standard. Il faut cependant garder à l'esprit qu'il n'est convergent que dans le très restrictif modèle de translation.

Une fois obtenu un estimateur convergent de  $V_{as}$ , l'inférence sur  $\beta_\tau$  est aisée. Un intervalle de confiance de niveau  $1 - \alpha$  sur  $\beta_\tau$  s'écrit ainsi :

$$IC_\alpha = \left[ \hat{\beta}_\tau - z_{1-\alpha/2} \sqrt{\hat{V}_{as}}, \hat{\beta}_\tau + z_{1-\alpha/2} \sqrt{\hat{V}_{as}} \right],$$

où  $z_{1-\alpha/2}$  est le quantile d'ordre  $1 - \alpha / 2$  d'une loi  $N(0, 1)$ . De même, la statistique de Wald  $T$  du test  $\beta_\tau = 0$  s'écrit  $T = n \hat{\beta}_\tau' \hat{V}_{as}^{-1} \hat{\beta}_\tau$ , avec  $T$  qui tend vers un  $\chi_p^2$  sous l'hypothèse nulle, où  $p$  est le nombre de variables explicatives.

#### 6 - Inférence par Bootstrap

Une autre possibilité pour faire de l'inférence est de recourir au bootstrap. Rappelons que le principe de bootstrap est de générer des échantillons « factices » par des tirages avec remise à partir de l'échantillon initial. Dans le cas du bootstrap standard, on applique l'algorithme suivant. De  $b = 1$  à  $B$  :

Tirer avec remise un échantillon de taille  $n$  à partir de l'échantillon initial  $(Y_i, X_i)_{i=1, \dots, n}$ . Soit  $(k_{b1}^*, \dots, k_{bn}^*)$  les indices correspondants aux observations tirées ;

$$\text{Calculer } \hat{\beta}_{tb}^* = \arg \min_{\beta} \sum_{j=1}^n \rho_\tau(Y_{k_{bj}^*} - X_{k_{bj}^*}' \beta).$$

On peut alors estimer la variance asymptotique par :

$$V_{as}^* = \frac{1}{B} \sum_{b=1}^B (\hat{\beta}_{tb}^* - \hat{\beta})^2.$$

Des intervalles de confiance ou tests peuvent être alors construits comme précédemment, en utilisant

l'approximation normale. Pour construire des intervalles de confiance, on peut également s'appuyer sur la *percentile bootstrap*. Soit  $q_u^*$  le quantile empirique d'ordre  $u$  de  $(\hat{\beta}_{\tau 1}^*, \dots, \hat{\beta}_{\tau B}^*)$ , on construit simplement l'intervalle de confiance par

$$IC_{1-\alpha} = [q_{\alpha/2}^*, q_{1-\alpha/2}^*].$$

Par rapport à l'estimateur (22), les méthodes de bootstrap ont l'avantage de ne pas supposer que le vrai modèle est un modèle de translation. Elles évitent également de devoir choisir le paramètre de lissage  $h_n$ , sachant que les résultats peuvent être sensibles à ce choix.

