

MESURER LES DIFFÉRENCES SALARIALES ENTRE DEUX GROUPES

Rapport méthodologique et paquet R “decr”

Journées suisses de la statistique, Zurich, 27 Août 2018

Sandro Petrillo (sandro.petrillo@ti.ch)
Oscar Gonzalez (oscar.gonzalez@ti.ch)

Repubblica e Cantone Ticino
Ufficio di statistica



Introduction

Quand on parle de salaires, on compare souvent les salaires de deux groupes d'individus, par exemple les hommes et les femmes ou suisses et étrangers. La comparaison directe entre les statistiques salariales de deux groupes peut cacher certains pièges. En effet, les travailleurs en comparaison peuvent être hétérogènes en termes de caractéristiques personnelles et de lieu de travail (profils de formation, postes hiérarchiques dans l'entreprise, branches économiques, etc.). Si ces différences ne sont pas suffisamment prises en compte, il existe un risque important d'interprétation erronée de l'ampleur et de la nature des disparités mesurées. À cet égard, nous avons développé un rapport méthodologique qui formalise une technique **non paramétrique** pour décomposer les différences observées entre les statistiques salariales de deux groupes d'individus en deux composantes:

- une attribuable au fait que les deux groupes ont des caractéristiques différentes (appelée **“partie expliquée”**);
- l'autre qui reflète les différences dans la structure salariale des deux groupes et n'est pas attribuable aux divergences dans la répartition des caractéristiques observées chez les travailleurs (appelée **“partie inexpliquée”**).

Caractéristiques principales de la méthode:

- Non paramétrique
- Utilisation des poids d'échantillonnage (pas présente dans d'autres méthodes)
- Applicable à toutes les statistiques qui dépendent de la distribution salariale
- Prise en compte explicite du **support commun**
- Paquet R “decr”

Étapes de l'analyse

- On observe deux groupes d'individus, A et B , dont on veut examiner les différences entre des statistiques salariales.
- Sur la base de certaines caractéristiques (par exemple la branche économique, la position hiérarchique, le niveau de formation, ...), on établit d'abord le **support commun**, c'est-à-dire les combinaisons de caractéristiques où l'on observe des individus des deux groupes.
- Dans le **support commun**, on estime un **facteur de repondération** pour chaque individu du groupe A ($\hat{\Psi}_A(\mathbf{x}_i) = \frac{\hat{f}_{X_B|S}(\mathbf{x}_i)}{\hat{f}_{X_A|S}(\mathbf{x}_i)}$).
- Les **facteurs de repondération** (multipliés par les **poids d'échantillonnage**) balancent la distribution des caractéristiques des individus du groupe A dans le support commun à celle des individus du groupe B . Les facteurs de repondération permettent donc d'estimer une **distribution des salaires “contre-factuelle”**, qui représente la distribution des salaires du groupe A comme s'ils avaient la même distribution des caractéristiques du groupe B .
- La distribution contre-factuelle des salaires permet l'estimation d'une **statistique salariale contre-factuelle** (moyenne, quantiles, ...), qui est la quantité fondamentale afin de **décomposer la différence des statistiques salariales des deux groupes**, dans le support commun, en une **partie expliquée** (par les différentes distributions des caractéristiques des deux groupes) et une **partie inexpliquée**.

Distribution contre-factuelle des salaires

La distribution contre-factuelle des salaires “mélange” la structure des salaires du groupe A avec les caractéristiques du groupe B .

$$F_{Y_A^C|S}(y) = \int F_{Y_A|X_A,S}(y|X) \cdot dF_{X_B|S}(X) \\ = \int F_{Y_A|X_A,S}(y|X) \cdot \Psi_A(X) \cdot dF_{X_A|S}(X),$$

où $\Psi_A(X) = dF_{X_B|S}(X)/dF_{X_A|S}(X)$ est un facteur de repondération.

Décomposition dans le support commun (S)

$$\Delta_{O|S}^v = \nu(F_{Y_A|D_A,S}) - \nu(F_{Y_B|D_B,S}) = \quad (1) \\ = (\nu(F_{Y_A|D_A,S}) - \nu(F_{Y_A^C|D_B,S})) + (\nu(F_{Y_A^C|D_B,S}) - \nu(F_{Y_B|D_B,S})) = \\ = \Delta_X^v + \Delta_S^v,$$

où Δ_X^v est la partie expliquée par les différentes distributions des caractéristiques des deux groupes dans le support commun et Δ_S^v est la partie inexpliquée.

Cas particulier des moyennes

$$\hat{\Delta}_O^\mu = \hat{Y}_A - \hat{Y}_B = \quad (2) \\ = \hat{p}_{S^C|A} \cdot \underbrace{(\hat{Y}_{A|S^C} - \hat{Y}_{A|S})}_{\hat{\Delta}_A^\mu \text{ (due aux individus hors support)}} + \underbrace{(\hat{Y}_{A|S} - \hat{Y}_{B|S})}_{\hat{\Delta}_{O|S}^\mu = \hat{\Delta}_X^\mu + \hat{\Delta}_S^\mu} + \hat{p}_{S^C|B} \cdot \underbrace{(\hat{Y}_{B|S} - \hat{Y}_{B|S^C})}_{\hat{\Delta}_B^\mu \text{ (due aux individus hors support)}}.$$

$$\hat{\Delta}_{O|S}^\mu = \underbrace{(\hat{Y}_{A|S} - \hat{Y}_{A|S^C})}_{\hat{\Delta}_X^\mu \text{ (expliquée)}} + \underbrace{(\hat{Y}_{A|S} - \hat{Y}_{B|S})}_{\hat{\Delta}_S^\mu \text{ (non expliquée)}}. \quad (3)$$

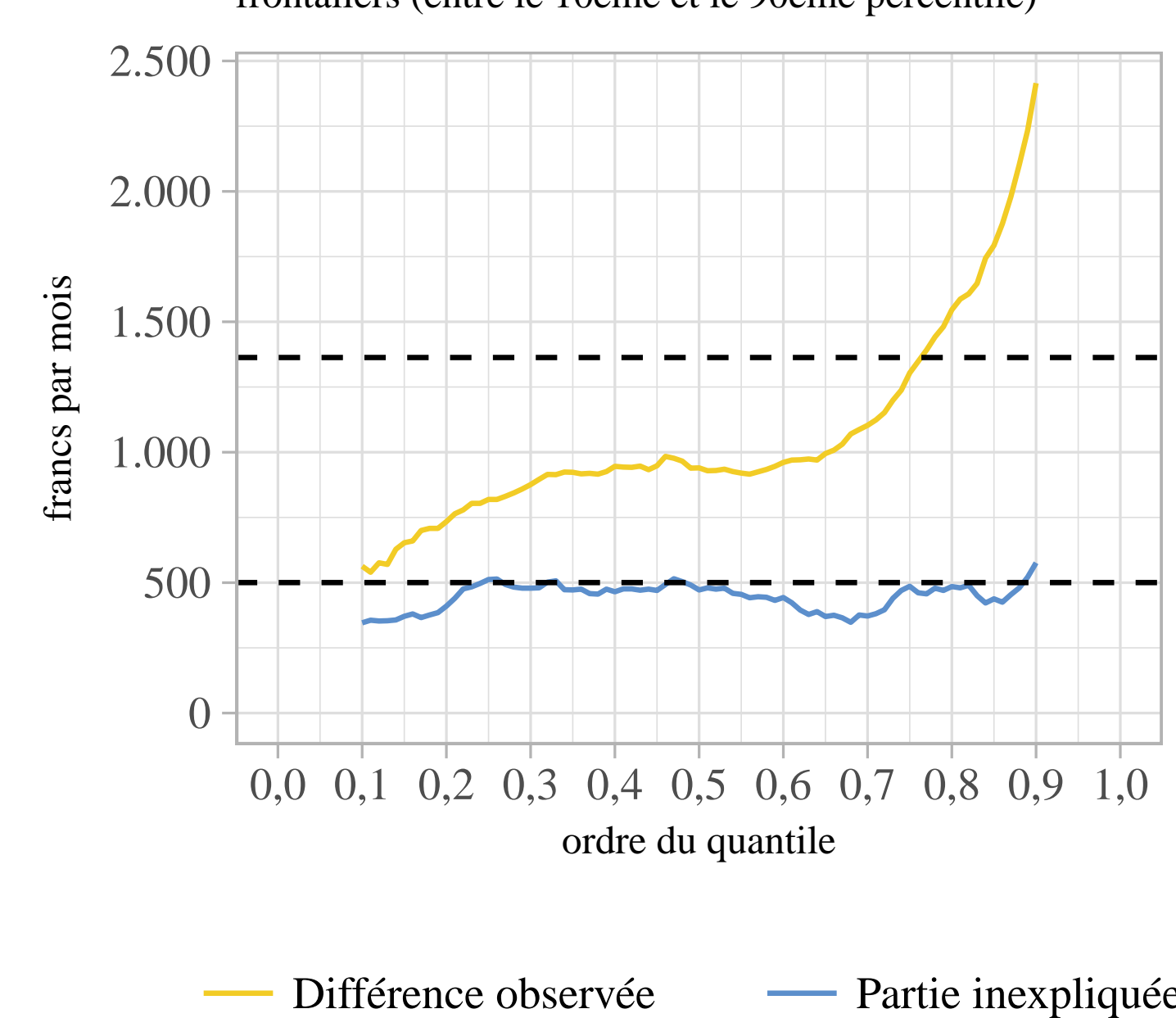
Application / Exemple

- En 2014, au Tessin, dans le secteur privé, différences entre résidents (A) et frontaliers (B) (Enquête suisse sur la structure des salaires, variable MBLS)
- Salaires moyens observés: Résidents: 6.249 francs/mois; Frontaliers: 4.886 francs/mois
- On observe donc une différence de 1.363 francs/mois
- Si on considère la section économique (17), la position dans la profession (4), le degré de formation (3/4), les classes d'âge de 10 ans (6) et les classes d'années de services dans la même entreprise (5), il y a 3.092 combinaisons (“profils”, “strates”) où on observe au moins un salarié.
- Dans 1.404 strates il existe au moins un individu de chacun des deux groupes (support commun). Les individus dans le support commun représentent le 87,2% des effectifs de l'échantillon.
- Décomposition de la différence des salaires moyens des résidents (A) et des frontaliers (B):

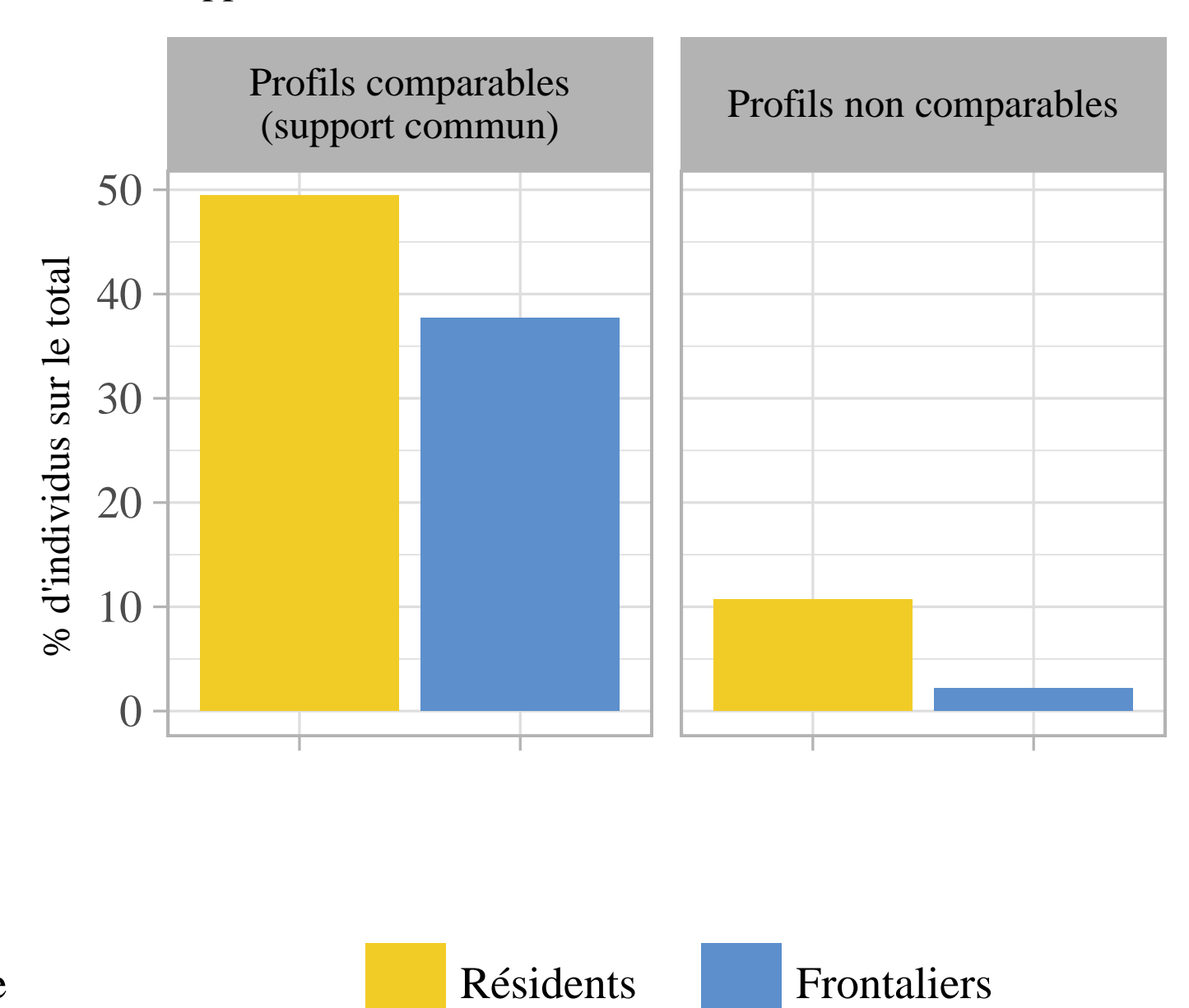
$$\begin{array}{r|rrrrr} \hat{\Delta}_O^\mu & \hat{\Delta}_A^\mu & \hat{\Delta}_X^\mu & \hat{\Delta}_S^\mu & \hat{\Delta}_B^\mu & \\ \hline 1.363 & 130 & 775 & 500 & -42 \end{array}$$

- La différence totale observée de 1.363 francs par mois ($\hat{\Delta}_O^\mu$) est décomposée en:
 - 130 francs ($\hat{\Delta}_A^\mu$) **expliquables** par le fait qu'il y a des résidents hors support qui ont un salaire moyen supérieur aux résidents du support commun,
 - 775 francs ($\hat{\Delta}_X^\mu$) **expliquables** par le fait que dans les strates comparables (support commun) les deux groupes se distribuent différemment,
 - 42 francs ($\hat{\Delta}_B^\mu$) **expliquables** par le fait qu'il y a des frontaliers hors support qui ont un salaire moyen supérieur aux frontaliers du support commun,
 - 500 francs ($\hat{\Delta}_S^\mu$) **inexpliquables** par les différentes caractéristiques considérées des deux groupes. En d'autres termes, cette composante est la différence “à parité de caractéristiques”.

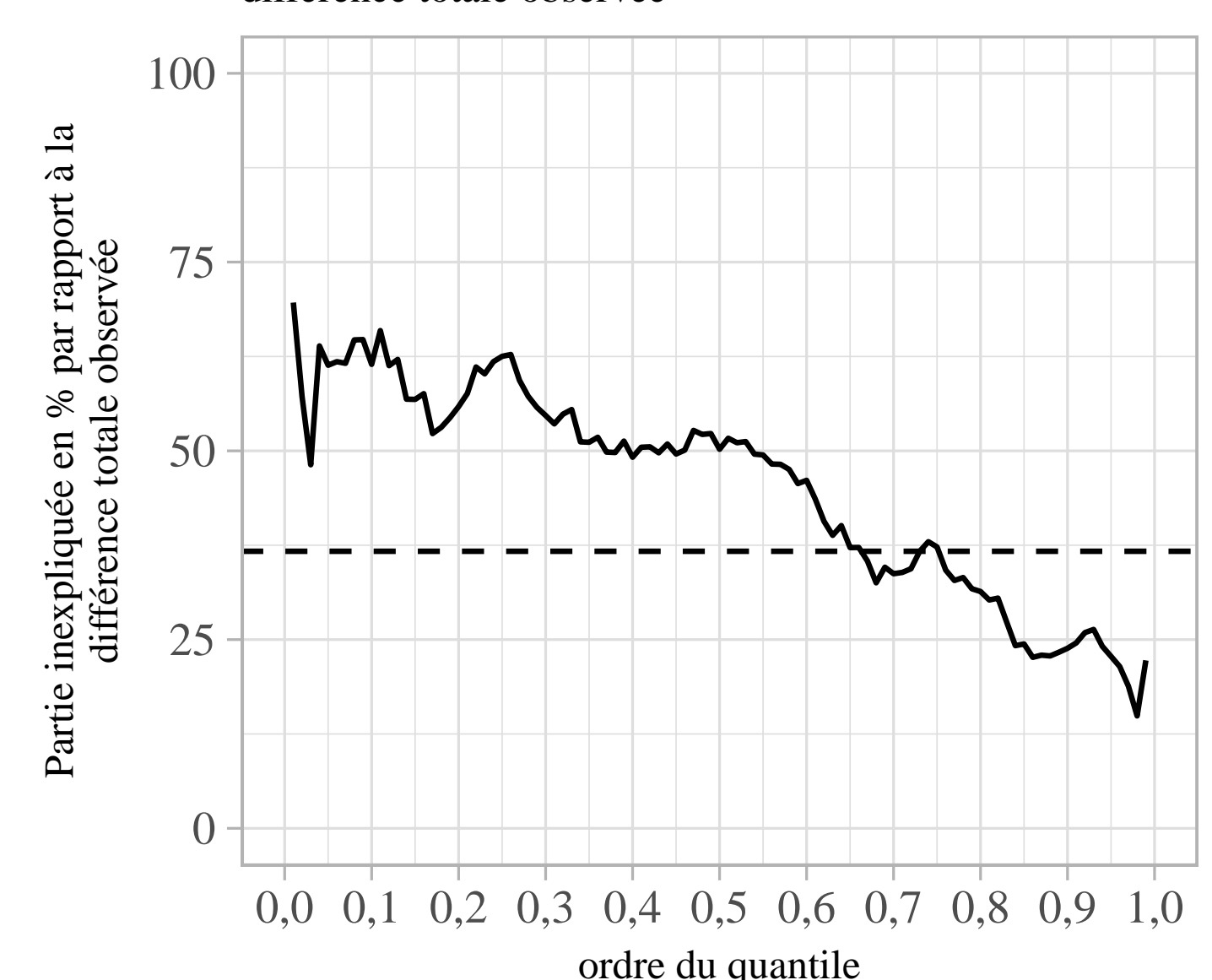
Différences observées et partie inexpliquée des différences entre les percentiles des salaires des résidents et des frontaliers (entre le 10ème et le 90ème percentile)



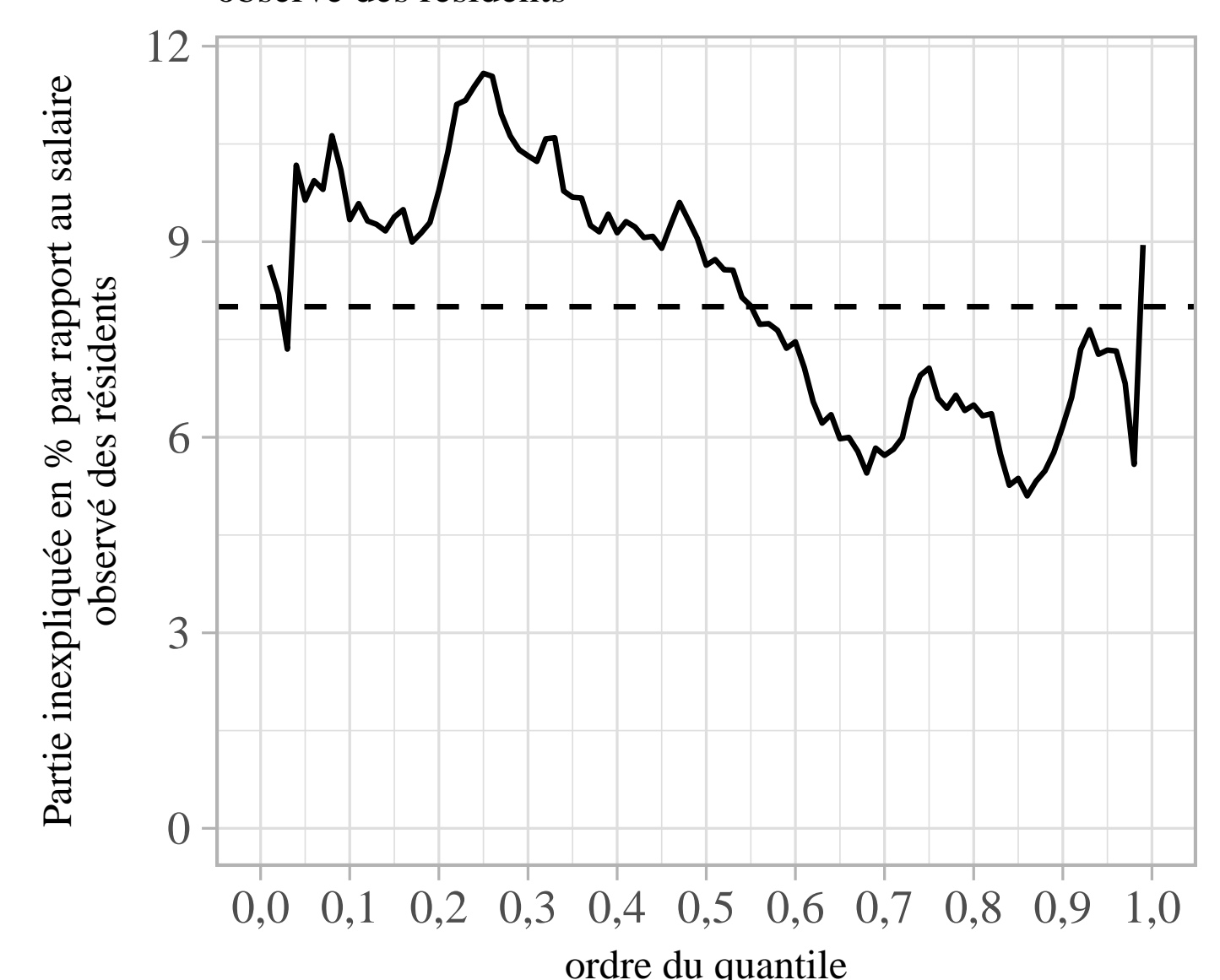
Pourcentage de résidents et frontaliers avec des profils comparables (support commun) et non comparables (hors support)



Partie inexpliquée en % par rapport à la différence totale observée



Partie inexpliquée en % par rapport au salaire observé des résidents



Exemple de code R avec le paquet “decr”

```
# install.packages("devtools")
devtools::install_github("gibonet/decr")
library(decr)
data(invented_wages)
str(invented_wages)

example(nopodec) # différence entre les salaires moyens
example(dec_)    # différence entre les quantiles d'un niveau
example(dec_all_) # plusieurs quantiles en une fois
```

Principales références bibliographiques

- DiNardo, John, Nicole M. Fortin, and Thomas Lemieux. 1996. “Labor Market Institutions and the Distribution of Wages, 1973-1992: A Semiparametric Approach.” *Econometrica* 64 (5):1001–44.
- Iacus, Stefano M., Gary King, and Giuseppe Porro. 2011. “Causal Inference Without Balance Checking: Coarsened Exact Matching.” *Political Analysis* 20:1–24.
- Nopo, Hugo. 2008. “Matching as a Tool to Decompose Wage Gaps.” *Review of Economics and Statistics* 90 (2):290–99.