
MÉTHODES ÉCONOMÉTRIQUES DE DÉCOMPOSITION DES INÉGALITÉS - DE LA THÉORIE À LA PRATIQUE

Sophie MAILLARD(), Béatrice BOUTCHENIK(*)(**)*

() INSEE – SSPLab*

*(**) Université Paris-Dauphine*

`sophie.maillard@insee.fr`

Keywords. Décomposition, Distribution contrefactuelle, Inégalités, Structure des salaires, Écart de salaire, Discrimination.

Résumé

Les méthodes de décomposition sont des outils standards pour l'analyse statistique des différences entre deux groupes, notamment salariales. Le modèle canonique d'Oaxaca-Blinder (Oaxaca 1973, Blinder 1973) propose ainsi une décomposition des écarts observés entre deux populations en une part expliquée par les caractéristiques observables de ces deux groupes, ou effet de composition, et en une part inexpliquée. Isoler cet écart inexpliqué permet de mettre en avant d'éventuels phénomènes discriminatoires, sous certaines hypothèses que nous nous attachons à clarifier. Plusieurs méthodes ont par ailleurs été proposées pour étendre le cadre classique d'Oaxaca-Blinder à une analyse plus complète des écarts, en particulier pour des variables catégorielles (Fairlie, 2005) et pour l'ensemble de la distribution de variables continues (Fortin, Firpo et Lemieux, 2011). Nous portons un intérêt particulier à cette extension aux distributions : celle-ci permet de mettre en évidence des effets hétérogènes, et notamment des mécanismes de "plafond de verre" ou de "plancher collant", suivant que les écarts se creusent dans le haut ou le bas de la distribution de la variable d'intérêt. Les méthodes correspondantes ont fait l'objet de nombreux développements récents (Chernozhukov et al., 2013 ; Firpo, Fortin et Lemieux, 2009). Nous discutons la mise en œuvre et la pertinence de ces méthodes et nous les illustrons à partir des données de l'Enquête Emploi en Continu (pour les années 2013 à 2016), pour l'exemple des disparités de salaire entre hommes et femmes et entre descendants d'immigrés et personnes sans ascendance migratoire.

Abstract

Decomposition methods are extensively used to analyze differentials between groups, typically in terms of wages. The seminal model of Oaxaca-Blinder (Oaxaca 1973, Blinder 1973) proposes a framework to decompose observed discrepancies between groups into an explained part which can be related to differences in observed characteristics, or composition effect, and an unexplained part. Isolating the unexplained gap can allow to identify discriminations under a number of assumptions that we explicit. Futhermore, different approaches extend the canonical model to categorical variables (Fairlie, 2005) and all along the distribution of a continuous variable. This extension of decomposition methods to distributions allows us to evidence heterogeneous effects, for instance “glass ceiling” or “sticky floor” phenomena, whether gaps widen at the top or at the bottom of the distribution of the variable of interest. The corresponding approaches have been the object of numerous recent developments (Fortin, Firpo and Lemieux, 2011 ; Chernozhukov et al., 2013 ; Firpo, Fortin and Lemieux, 2009). We discuss the application and the relevance of these methods and illustrate them using French labor force survey data between 2013 and 2016, with a case study on the gender and ethnic wage gap.

Introduction

On appelle méthodes de décomposition l'ensemble des techniques visant à séparer une différence - par exemple de revenus - entre deux groupes en une part liée à des caractéristiques observées individuelles différentes, la *part expliquée* ou *effet de composition*, et une part résiduelle à caractéristiques observables égales, la *part inexpliquée*. Les exemples canoniques de décomposition étudient les écarts de salaires moyens entre hommes et femmes (Oaxaca, 1973) et entre individus “blancs” et “noirs” aux États-Unis (Blinder, 1973) : une partie de ces écarts pourrait par exemple être attribuée à des différences dans les niveaux d'éducation, ou d'expérience sur le marché du travail. Les méthodes de décomposition se sont imposées comme des outils essentiels dans l'étude des inégalités, du fait notamment de leur facilité de mise en œuvre. Celles-ci permettent de mesurer l'ampleur des effets de composition dans l'inégalité, et d'analyser de façon détaillée les contributions de plusieurs facteurs à l'écart total de revenus, et de porter ainsi un diagnostic fin sur les mécanismes de formation des inégalités. De nombreuses applications sortant du cadre de l'étude des discriminations se prêtent ainsi à la mise en œuvre de ces méthodes, selon le type de groupes que l'on souhaite comparer. On pourra par exemple décomposer l'écart entre le revenu moyen dans un département donné et le revenu moyen au niveau national, afin de comprendre les situations inégales de différents territoires (Bertran, 2017 pour un exemple sur les revenus d'activité des non-salariés). Dans le cadre d'une comparaison internationale, on sera amené à analyser l'écart pour une même mesure entre différents pays, contrastés deux à deux. Les “groupes” étudiés peuvent également être deux périodes distinctes, en lesquelles on considère une même grandeur, cherchant à comprendre les déterminants de son évolution : Audenaert et al., 2014 décomposent par exemple la croissance du salaire moyen en France dans les années 2000, isolant ce qui relève des évolutions de composition de la population salariée. Les méthodes de décomposition sont souvent liées à l'analyse des discriminations. Pour autant, le cadre dans lequel ce type d'interprétation est possible nécessite des hypothèses beaucoup plus

fortes que celles requises lorsqu'on souhaite seulement isoler un effet de composition.

Les méthodes de décomposition ont connu ces dernières années un fort regain d'intérêt, notamment dans un contexte de hausse rapide des inégalités de salaire aux États-Unis (Firpo et al., 2007). Premièrement, une importante réflexion a eu lieu sur les conditions de l'identification d'une discrimination, notamment par l'analogie avec la notion d'effet de traitement, empruntée aux travaux d'évaluation de politiques publiques. Mesurer une discrimination entre deux groupes en contrôlant de leurs caractéristiques peut s'apparenter à estimer l'effet d'un traitement en comparant un groupe traité et un groupe de contrôle. La littérature sur les décompositions s'est ainsi demandé sous quelles conditions il était possible d'interpréter la part inexpliquée de la décomposition comme l'effet pur de l'appartenance au groupe considéré, et *in fine* comme une mesure de discrimination. Deuxièmement, la réflexion sur les méthodes de décomposition a permis d'élargir la palette des outils hors du cadre initialement proposé par Blinder et Oaxaca, s'appuyant sur la régression linéaire et permettant la décomposition de l'écart entre moyennes par groupes pour une variable continue. Des méthodes ont ainsi été développées permettant de s'intéresser non plus à une variable continue mais à une variable dichotomique, d'une part ; et d'analyser l'ensemble de la distribution des variables continues, d'autre part.

Fortin, Lemieux, and Firpo (2011) ont détaillé dans un article de référence ces nouvelles méthodes, et plus généralement le cadre théorique entourant les méthodes de décomposition¹. Nous nous appuyons ici sur ce document pour en proposer une traduction pratique, tout en insistant sur les questions qu'il est nécessaire de se poser afin d'interpréter correctement les résultats issus des méthodes de décomposition. Le cadre théorique correspondant, ainsi que certains approfondissements, sont renvoyés en encadré. Des codes R sont proposés pour les méthodes les plus facilement implémentables.

Lorsque la variable d'intérêt est continue (salaires, revenus, patrimoine, heures travaillées, notes à un examen...), on considère le plus souvent sa moyenne par groupe. La décomposition de l'écart entre moyennes, qui correspond au cadre classique d'Oaxaca-Blinder, est abordée dans la section 1. On présente en section 2 les précautions à prendre dans l'interprétation des résultats, en particulier les conditions sous lesquelles un écart inexpliqué peut être interprété comme un effet causal de l'appartenance à un groupe plutôt qu'à l'autre. La variable d'intérêt peut également être une variable dichotomique : le fait d'être au chômage, d'être actif ou encore d'avoir un emploi stable. On cherchera alors par exemple à comprendre l'écart entre les taux de chômage mesurés dans un groupe et dans l'autre. Ces cas sont traités dans la section 3. Lorsque la variable est continue, on s'intéresse parfois à d'autres statistiques que sa moyenne, et notamment aux écarts existant en différents points de la distribution : le lecteur pourra se référer à la section 4 pour les différentes méthodes de décomposition des écarts entre distributions.

1. Pour que le lecteur puisse s'y reporter plus facilement, nous conservons des notations proches de celles utilisées par Fortin et al. (2011).

1 La décomposition de l'écart de moyennes entre deux groupes

Dans cette partie, nous présentons la méthode la plus classique de décomposition des inégalités à la moyenne, dans le cas où la variable d'intérêt est continue : la décomposition dite d'*Oaxaca-Blinder*.

1.1 Le modèle classique d'Oaxaca-Blinder

On considère ici une variable continue Y , dont on observe un ensemble de K déterminants individuels X_1, X_2, \dots, X_K . On souhaite étudier l'écart entre les moyennes de Y selon deux groupes A et B , en lien avec le fait que ces deux groupes présentent des caractéristiques observables différentes. Par exemple, la variable Y pourrait correspondre au salaire, les variables X au niveau d'éducation, à l'expérience sur le marché du travail, etc., et les groupes A et B aux hommes et aux femmes. On modélise séparément, dans le groupe A et le groupe B , une relation linéaire entre la variable Y et ses déterminants :

$$Y_i = \beta_{A0} + \sum_{k=1}^K X_{ik} \beta_{Ak} + v_{iA}, \forall i \in A$$

$$Y_i = \beta_{B0} + \sum_{k=1}^K X_{ik} \beta_{Bk} + v_{iB}, \forall i \in B$$

Une fois les paramètres de chacun des deux modèles estimés, on peut alors écrire, en notant $\overline{Y_B}$ et $\overline{Y_A}$ le salaire moyen dans chaque groupe :

$$\overline{Y_A} = \hat{\beta}_{A0} + \sum_{k=1}^K \overline{X_{Ak}} \hat{\beta}_{Ak}$$

$$\overline{Y_B} = \hat{\beta}_{B0} + \sum_{k=1}^K \overline{X_{Bk}} \hat{\beta}_{Bk}$$

Le salaire moyen peut différer d'un groupe à l'autre pour deux raisons : d'une part, parce que les *caractéristiques* moyennes ne sont pas les mêmes dans le groupe A et le groupe B ; d'autre part, parce que les *valorisations* de ces caractéristiques (les $(\hat{\beta}_{g,k})_{k=1 \dots K}, g = A, B$), ainsi que les constantes des deux modèles, sont différentes. On pourra ainsi décomposer l'écart entre $\overline{Y_B}$ et $\overline{Y_A}$ de la façon suivante :

$$\begin{aligned} \overline{Y_B} - \overline{Y_A} &= \hat{\beta}_{B0} + \sum_{k=1}^K \overline{X_{Bk}} \hat{\beta}_{Bk} - \hat{\beta}_{A0} - \sum_{k=1}^K \overline{X_{Ak}} \hat{\beta}_{Ak} \\ &= \underbrace{\sum_{k=1}^K (\overline{X_{Bk}} - \overline{X_{Ak}}) \hat{\beta}_{Bk}}_{\hat{\Delta}_X \text{ (expliqué)}} + \underbrace{(\hat{\beta}_{B0} - \hat{\beta}_{A0}) + \sum_{k=1}^K \overline{X_{Ak}} (\hat{\beta}_{Bk} - \hat{\beta}_{Ak})}_{\hat{\Delta}_S \text{ (inexpliqué)}} \end{aligned} \quad (1)$$

$\hat{\Delta}_X$ renvoie à la partie de l'écart de salaire liée à l'écart de caractéristiques observables entre les deux groupes, écart que l'on valorise ici selon les paramètres estimés pour le groupe B : on appellera cette grandeur l'écart expliqué (ou *effet de composition*). $\hat{\Delta}_S$ correspond à la part liée à l'écart de valorisation de caractéristiques (et à l'écart de constante), valorisations qui ici sont appliquées aux caractéristiques du groupe A . On désignera ce terme comme "écart inexpliqué".

Encadré 1 : L'ambiguïté du terme d'*effet de structure*

La littérature des méthodes de décomposition appelle aussi l'écart inexpliqué $\hat{\Delta}_S$ le "wage structure effect". Cette appellation vient de l'hypothèse qu'il existe une fonction structurelle des salaires qui diffère entre les groupes comparés. Autrement dit, la structure à laquelle il est fait référence dans ce terme est une structure de *valorisation des caractéristiques*, et non une structure de *caractéristiques*- comme l'entend le plus souvent le langage courant dans le terme d'"effet de structure". Cela pouvant être source de confusion, on évitera ici de parler d'effet structurel.

La notion de structure est aussi centrale dans d'autres méthodes de décomposition, comme les approches dites structurelles géographiques. Celles-ci permettent d'analyser les différences d'évolution entre territoires, entre ce qui tient des *structures* sectorielles spécifiques et des effets résiduels (compétitivité locale, capacités d'innovation, etc) qui correspondent à la différence entre la croissance de l'ensemble des territoires et de la zone d'intérêt à *structure* productive donnée. Pour plus d'éléments sur ces méthodes, on pourra se reporter à Kubrak (2018).

1.2 Application : la décomposition agrégée

On illustre la décomposition d'Oaxaca-Blinder par l'étude des différences de salaires entre hommes et femmes, à partir de l'enquête Emploi entre 2013 et 2016. Si l'on se réfère à la décomposition (1), et que l'on souhaite décomposer $\overline{Y}_B - \overline{Y}_A$, B correspondra aux hommes et A aux femmes. La variable d'intérêt est le logarithme du salaire mensuel net. Celui-ci vaut en moyenne 7.572 chez les hommes, et 7.273 chez les femmes, soit un écart de 0.299.

On introduit comme variables explicatives l'expérience potentielle et son carré, le niveau d'études en 6 postes (Diplôme supérieur à baccalauréat + 2 ans, Baccalauréat + 2 ans, Baccalauréat ou brevet professionnel ou autre diplôme de ce niveau, CAP, BEP ou autre diplôme de ce niveau, Brevet des collèges, Aucun diplôme ou certificat d'études primaires), le secteur d'activité (référence=secteur du commerce et de l'hébergement-restauration), une indicatrice d'être à temps partiel et l'ancienneté dans l'entreprise en 4 modalités (référence= moins d'un an). Afin d'effectuer la décomposition d'Oaxaca-Blinder correspondante, on estime les coefficients de l'équation de salaire dans chacun des groupes.

On s'intéresse dans un premier temps au partage *global* entre effet de composition $\hat{\Delta}_X$ et écart inexpliqué $\hat{\Delta}_S$: c'est la décomposition agrégée. Pour ce faire, on a en fait seulement besoin

d'estimer le modèle chez les hommes², mais la comparaison étant d'intérêt on procède aussi à l'estimation chez les femmes.

```
modele.A <- lm(logsal ~ exp_mtra + exp_mtra2resc +
               as.factor(ddipl) + tpartiel + secteur0Q +
               secteurBE + secteurRU + secteurFZ + secteurMN +
               secteurAZ + secteurKZ + secteurJZ + secteurLZ +
               ancentr44 + ancentr43 + ancentr42,
               data = data[data$sex==1,])
coeffs.A <- modele.A$coefficients

modele.B <- lm(logsal ~ exp_mtra + exp_mtra2resc +
               as.factor(ddipl) + tpartiel + secteur0Q +
               secteurBE + secteurRU + secteurFZ + secteurMN +
               secteurAZ + secteurKZ + secteurJZ + secteurLZ +
               ancentr44 + ancentr43 + ancentr42,
               data = data[data$sex==0,])
coeffs.B <- modele.B$coefficients

round(cbind(coeffs.A, coeffs.B), 3)

##                coeffs.A coeffs.B
## (Intercept)          6.721    6.827
## exp_mtra              0.010    0.025
## exp_mtra2resc        -0.016   -0.035
## as.factor(ddipl)1     0.684    0.718
## as.factor(ddipl)3     0.504    0.434
## as.factor(ddipl)4     0.314    0.287
## as.factor(ddipl)5     0.194    0.134
## as.factor(ddipl)6     0.175    0.169
## tpartiel             -0.510   -0.672
## secteur0Q            -0.044   -0.101
## secteurBE             0.091    0.060
## secteurRU            -0.245   -0.154
## secteurFZ             0.030    0.036
## secteurMN             0.008    0.001
## secteurAZ            -0.151   -0.119
## secteurKZ             0.136    0.183
## secteurJZ             0.197    0.124
## secteurLZ             0.011   -0.060
## ancentr44             0.367    0.202
## ancentr43             0.209    0.119
```

2. cf *infra* pour plus de détails.

```
## ancentr42          0.113    0.083
```

On calcule ensuite les moyennes pour chaque variable, pour chacun des deux groupes. Dans le cas des variables catégorielles, ici le diplôme, on a besoin des proportions pour chacune des modalités (hors référence). On réécrit également les variables catégorielles comme autant d'indicateurs qu'il y a de modalités, car cela simplifie le calcul des écarts expliqué et inexpliqué. Pour faire cette transformation automatiquement, on peut utiliser la fonction `model.matrix`.

```
X.A <- model.matrix(~ exp_mtra + exp_mtra2resc
+ as.factor(ddipl) + tpartiel + secteur0Q
+ secteurBE + secteurRU + secteurFZ + secteurMN
+ secteurAZ + secteurKZ + secteurJZ + secteurLZ
+ ancentr44 + ancentr43 + ancentr42,
data = data[data$sex==1,])
```

#on applique la fonction moyenne pour chaque variable

```
X.moy.A<-apply(X.A,2,mean)
```

```
X.B <- model.matrix(~ exp_mtra + exp_mtra2resc
+ as.factor(ddipl) + tpartiel + secteur0Q
+ secteurBE + secteurRU + secteurFZ + secteurMN
+ secteurAZ + secteurKZ + secteurJZ + secteurLZ
+ ancentr44 + ancentr43 + ancentr42,
data = data[data$sex==0,])
```

```
X.moy.B<-apply(X.B,2,mean)
```

```
round(cbind(X.moy.A,X.moy.B),3)
```

```
##                X.moy.A X.moy.B
## (Intercept)      1.000  1.000
## exp_mtra         22.498  22.230
## exp_mtra2resc    6.413  6.273
## as.factor(ddipl)1 0.218  0.192
## as.factor(ddipl)3 0.186  0.141
## as.factor(ddipl)4 0.201  0.185
## as.factor(ddipl)5 0.232  0.304
## as.factor(ddipl)6 0.055  0.049
## tpartiel         0.306  0.056
## secteur0Q        0.482  0.211
## secteurBE        0.086  0.227
## secteurRU        0.065  0.030
## secteurFZ        0.013  0.100
## secteurMN        0.086  0.089
```

```
## secteurAZ          0.007  0.016
## secteurKZ          0.041  0.028
## secteurJZ          0.017  0.037
## secteurLZ          0.014  0.011
## ancentr44          0.505  0.510
## ancentr43          0.172  0.171
## ancentr42          0.220  0.219
```

Pour retrouver l'effet de composition défini en (1), il reste seulement à appliquer les coefficients estimés chez les hommes aux différences entre les caractéristiques moyennes chez les hommes et chez les femmes et à sommer pour toutes les variables. Cela donne :

```
sum((X.moy.B- X.moy.A)*coeffs.B)

## [1] 0.177
```

à rapporter à un écart total de log salaire de 0.299 entre hommes et femmes. L'effet de composition représente ainsi 59.1 % de l'écart total de salaire observé entre les sexes. Autrement dit, 59.1 % de l'écart de salaire observé entre hommes et femmes à partir de l'enquête Emploi peut être attribué à des caractéristiques moyennes différentes entre les sexes. On peut vérifier que, mécaniquement, l'écart inexplicé correspond bien à 0.122 :

```
sum(X.moy.A*(coeffs.B-coeffs.A))

## [1] 0.122
```

On peut voir que dans l'exemple précédent, où l'on s'intéresse uniquement à la décomposition *agrégée*, il suffit en fait d'estimer le jeu de coefficients $(\hat{\beta}_{B,k})_{k=1\dots K}$ des hommes pour obtenir la décomposition souhaitée. En effet, on peut réécrire :

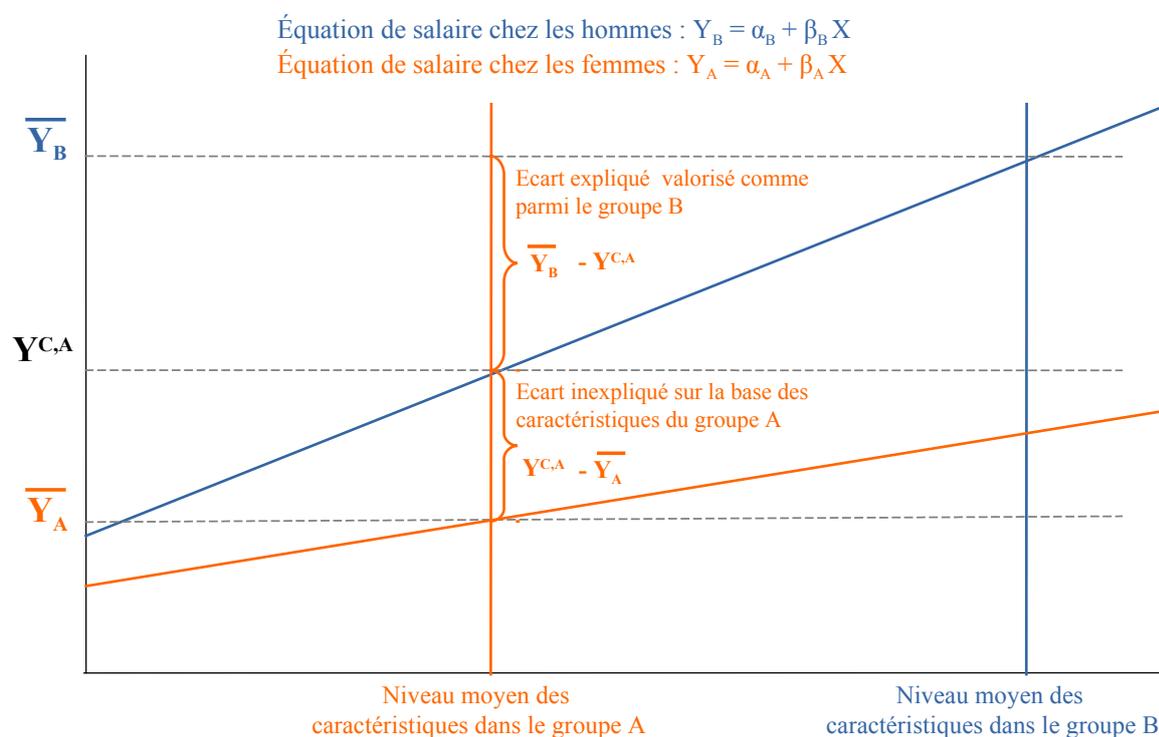
$$\overline{Y}_B - \overline{Y}_A = \underbrace{\overline{Y}_B - \sum_{k=0}^K \overline{X}_{Ak} \hat{\beta}_{Bk}}_{\hat{\Delta}_X} + \underbrace{\sum_{k=0}^K \overline{X}_{Ak} \hat{\beta}_{Bk} - \overline{Y}_A}_{\hat{\Delta}_S} \quad (2)$$

On ne s'appuie ici que sur les $\hat{\beta}_{Bk}$, et non sur les $\hat{\beta}_{Ak}$: cette formulation de la décomposition "agrégée" est utile lorsque l'un des deux groupes considérés comporte des effectifs très faibles, ce qui conduirait à une faible précision si l'on devait s'appuyer sur les coefficients estimés dans ce groupe. Cette remarque n'est pas toujours vraie, que l'on veuille aller plus loin que la décomposition agrégée ou que l'on souhaite considérer une autre valorisation de référence des caractéristiques.

1.3 Références de la décomposition

Dans la formule (1), on a implicitement introduit un salaire “contrefactuel”³ $Y^{C,A}$ valant $\hat{\beta}_{B0} + \sum_{k=1}^K \overline{X_{Ak}} \hat{\beta}_{Bk}$. Il correspond au salaire obtenu pour les caractéristiques observables moyennes du groupe A valorisées comme dans le groupe B . La question posée par ce contrefactuel peut se formuler ainsi : que gagneraient les individus du groupe A si leurs caractéristiques étaient valorisées de la même manière que pour les B ? L'écart entre ce terme et le salaire moyen du groupe B , $\hat{\beta}_{B0} + \sum_{k=1}^K \overline{X_{Bk}} \hat{\beta}_{Bk}$, résulte uniquement de différences de caractéristiques : on retrouve l'effet de composition. L'écart entre $Y^{C,A}$ et le salaire moyen du groupe A correspond à l'écart inexpliqué.

Figure 1 – Décomposition d'écart moyen de salaire entre les groupes B et A



Ceci est illustré sur la figure 1 qui présente un cas simple où l'on dispose d'une seule variable observable X . Les accolades en orange présentent d'une part l'écart entre contrefactuel $Y^{C,A}$ et salaire moyen du groupe B (hommes) (écart expliqué), d'autre part l'écart entre salaire moyen des A (femmes) et contrefactuel $Y^{C,A}$ (écart inexpliqué).

Un contrefactuel alternatif à $Y^{C,A}$ correspondrait au salaire qu'aurait le groupe B si ses caractéristiques étaient valorisées comme celles du groupe A , c'est-à-dire $\hat{\beta}_{A0} + \sum_{k=1}^K \overline{X_{Bk}} \hat{\beta}_{Ak}$. On note ce contrefactuel $Y^{C,B}$ et on dessine en bleu les accolades illustrant la décomposition suivant ce contrefactuel sur la figure 1. La décomposition correspondante est la suivante :

3. On emploie ici le terme contrefactuel à la façon de Fortin et al. (2011) pour désigner le salaire de référence de la décomposition- celui qu'aurait par exemple les femmes si, à caractéristiques observables inchangées, celles-ci étaient valorisées comme parmi les hommes. Ce salaire de référence ne s'interprète pas de façon causal.

$$\overline{Y_B} - \overline{Y_A} = \underbrace{\sum_{k=1}^K (\overline{X_{Bk}} - \overline{X_{Ak}}) \hat{\beta}_{Ak}}_{\hat{\Delta}_X \text{ (expliqué)}} + \underbrace{(\hat{\beta}_{B0} - \hat{\beta}_{A0}) + \sum_{k=1}^K \overline{X_{Bk}} (\hat{\beta}_{Bk} - \hat{\beta}_{Ak})}_{\hat{\Delta}_S \text{ (inexpliqué)}} \quad (3)$$

Ici, l'écart de caractéristiques entre les deux groupes est donc valorisé selon les coefficients β_A , et non selon les β_B comme c'était le cas dans la décomposition 1. On peut toutefois remarquer que rien n'empêche de considérer n'importe quel autre vecteur de coefficients β_Ω comme la référence de la décomposition. On pourra par exemple choisir comme coefficients β_Ω ceux estimés sur l'ensemble de la population. L'écart inexpliqué comprend alors un terme supplémentaire, en effet dans ce cas la décomposition s'écrit :

$$\begin{aligned} \overline{Y_B} - \overline{Y_A} = & \underbrace{(\hat{\beta}_{B0} - \hat{\beta}_{A0}) + \sum_{k=1}^K \overline{X_{Bk}} (\hat{\beta}_{Bk} - \hat{\beta}_{\Omega k}) + \sum_{k=1}^K \overline{X_{Ak}} (\hat{\beta}_{\Omega k} - \hat{\beta}_{Ak})}_{\hat{\Delta}_S^v} \\ & + \underbrace{\sum_{k=1}^K (\overline{X_{Bk}} - \overline{X_{Ak}}) \hat{\beta}_{\Omega k}}_{\hat{\Delta}_X^v} \end{aligned} \quad (4)$$

On détaille en section 2.2 les questions à se poser pour bien choisir la référence de la décomposition.

Cette question du salaire de référence permet de faire le lien entre méthodes de décomposition et une autre méthode courante d'analyse des écarts de salaire entre deux groupes consistant à introduire simplement dans l'équation de salaire une indicatrice d'appartenance à l'un ou l'autre des groupes :

$$Y_i = \beta_0 + \sum_{k=1}^K X_{ik} \beta_k + \mathbf{1}_{i \in B} \beta_B$$

Cette méthode permet, en contrôlant des différences de caractéristiques observables entre les groupes, d'obtenir une estimation alternative de l'écart inexpliqué, correspondant à $\hat{\beta}_B$. On peut voir que cet écart inexpliqué peut être retrouvé en utilisant comme référence dans une décomposition d'Oaxaca-Blinder une valorisation de référence commune entre les deux groupes sauf pour la constante. Autrement dit, la méthode de l'indicatrice est un cas particulier de la méthode d'Oaxaca-Blinder.

1.4 La décomposition détaillée de l'effet de composition

Afin d'avoir une vision plus fine des mécanismes jouant sur l'effet de composition, il est possible de détailler celui-ci variable par variable. Ainsi, on peut considérer un à un au sein de $\hat{\Delta}_X$, chacun des termes liés à une variable explicative X_k en particulier :

$$\hat{\Delta}_X = \sum_{k=1}^K \hat{\Delta}_{X_k}$$

où pour chaque covariable X_k , $\hat{\Delta}_{X_k}$ désigne sa contribution à l'écart expliqué

$$\hat{\Delta}_{X_k}^\nu = (\overline{X_{Bk}} - \overline{X_{Ak}}) \hat{\beta}_{Bk}.$$

Comme dans le cas simple de la décomposition agrégée, on n'a besoin d'estimer que les valorisations des caractéristiques du groupe B pour calculer chacun des termes de l'effet de composition.

On utilise à présent le package `Oaxaca` qui permet d'automatiser les calculs des écarts expliqué et inexpliqué, de comparer différentes références et de détailler l'analyse variable par variable. On pourra se reporter à Hlavac (2014) pour plus de détails. L'exemple d'application est le même que précédemment.

```
library("oaxaca")
```

On utilise la fonction `Oaxaca` pour renseigner le modèle linéaire sur lequel est fondé la décomposition et la variable permettant de distinguer les deux groupes à comparer. Par défaut, les erreurs sont calculées par bootstrap, à partir de 100 réplifications. On peut modifier ce paramètre en spécifiant le paramètre `R`.

```
results <- oaxaca(formula = logsal ~ exp_mtra + exp_mtra2resc
  + ddipl6 + ddipl5 + ddipl4 + ddipl3 + ddipl1
  + tpartiel + secteurOQ + secteurBE + secteurRU
  + secteurFZ + secteurMN + secteurAZ + secteurKZ
  + secteurJZ + secteurLZ + ancentr44 + ancentr43
  + ancentr42 | sex , data = data, R=50)
```

Une fois les paramètres de la décomposition estimés, on peut afficher différentes sorties, comme la composante `n` qui renvoie le nombre d'observations dans les deux groupes ou `y` qui donne les salaires moyens dans chaque groupe et la différence entre les deux. Plus intéressant, on peut afficher les résultats de la décomposition agrégée :

```
round(results$twofold$overall[,1:5], 3)
```

##	group.weight	coef(explained)	se(explained)	coef(unexplained)	se(unexplained)
## [1,]	0.000	0.133	0.002	0.166	0.002
## [2,]	1.000	0.177	0.002	0.122	0.003
## [3,]	0.500	0.155	0.002	0.144	0.002
## [4,]	0.489	0.154	0.002	0.144	0.002
## [5,]	-1.000	0.176	0.002	0.123	0.001
## [6,]	-2.000	0.144	0.002	0.155	0.002

La colonne `group.weight` indique à partir de quelle référence est calculée la décomposition :

- pour la ligne 0, les coefficients de référence sont estimés dans le groupe tel que la variable renseignée pour distinguer les deux populations (ici "sex") soit égale à 0. En l'occurrence, il s'agit des femmes. L'effet de composition obtenu avec cette décomposition correspond

à l'écart de salaire entre hommes et femmes lié à leur différence de caractéristiques lorsqu'on les valorisent comme chez les femmes.

- pour la ligne 1, on a les résultats de la décomposition avec les coefficients de référence estimés chez les hommes (variable "sex"=1). On peut remarquer qu'on retrouve exactement le même résultat que précédemment.
- 0.5 : moyenne (non pondérée) des coefficients estimés séparément dans chacun des groupes. Implicitement, cela revient à se rapporter à une valorisation moyenne des caractéristiques (Reimers, 1983).
- 0.489 : moyenne pondérée des coefficients estimés dans chaque groupe (Cotton, 1988). En fait, une telle définition voudrait plutôt que la pondération vaille alors de 0.511... Dans la version du package que nous utilisons, cette erreur semble demeurer.
- -1 : coefficients de référence estimés sur l'ensemble de la population, sans indicatrice de groupe (Neumark, 1988). Cette décomposition revient à considérer comme valorisation de référence une valorisation strictement identique des caractéristiques entre les deux groupes.
- -2 : coefficients de référence estimés sur l'ensemble de la population mais avec indicatrice de groupe (Jann, 2008). La référence considérée autorise donc seulement la constante du modèle à différer entre les deux groupes comparés.

Ainsi, pour le modèle `group.weight= 1`, on retrouve comme précédemment que l'écart expliqué vaut 0.177 et correspond à la différence entre le salaire que toucherait les femmes si elles avaient les caractéristiques moyennes des hommes valorisés comme chez les femmes et le salaire moyen effectivement observé chez les femmes. L'écart inexpliqué, de 0.122 point, correspond à la différence entre le salaire moyen des hommes et le salaire que toucherait les femmes si elles avaient les caractéristiques moyennes des hommes valorisés comme chez les femmes. Cette répartition entre expliqué et inexpliqué varie avec la référence retenue. Par exemple, on trouve, en considérant comme *contrefactuel* de la décomposition le salaire que toucheraient les femmes si leurs caractéristiques étaient valorisées comme celles des hommes (`group.weight= 0`), un effet de composition de 0.133 et un écart inexpliqué de 0.166.

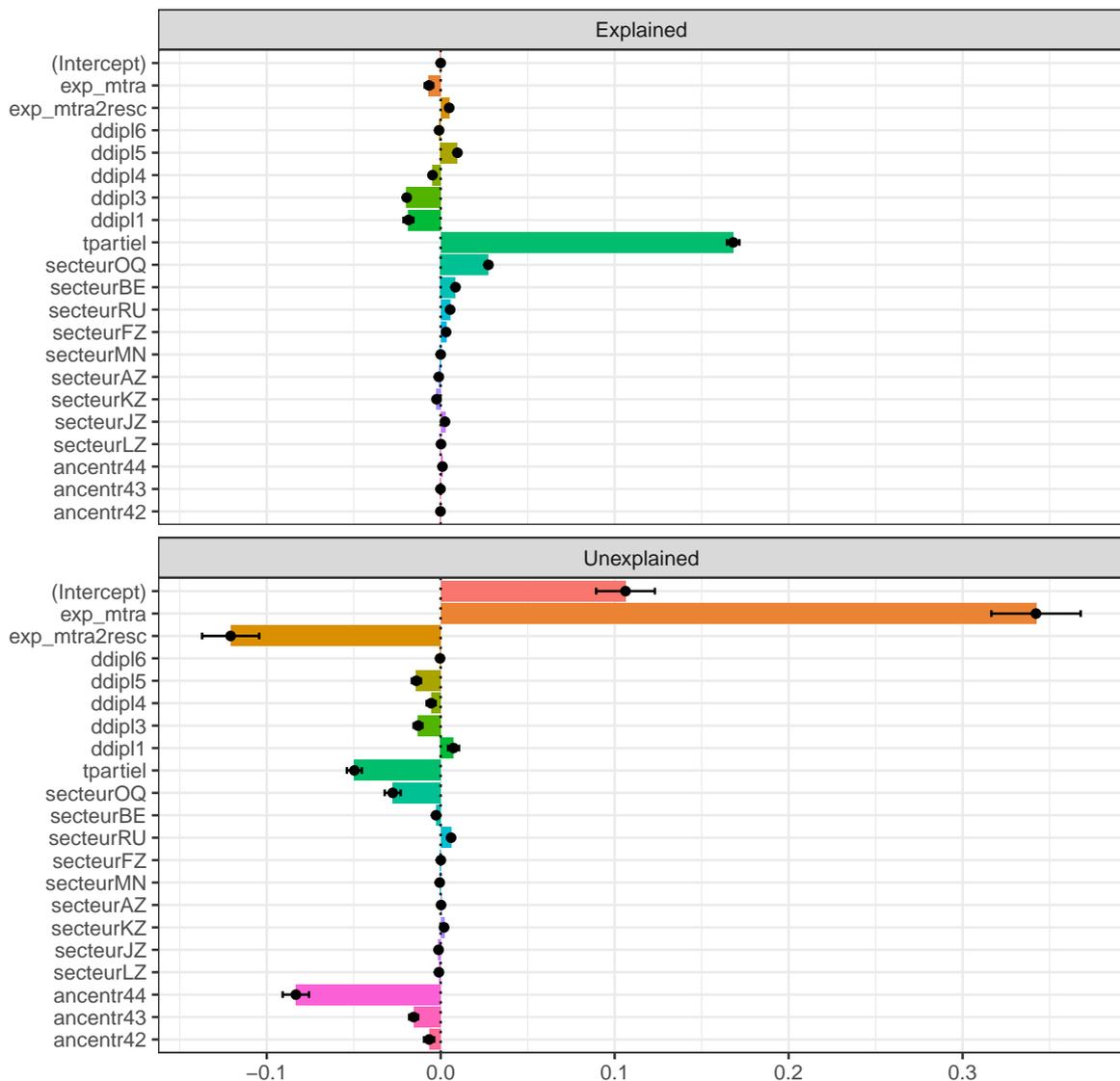
Pour chaque ensemble de coefficients de référence on peut ensuite afficher les résultats détaillés, c'est-à-dire la contribution de chaque variable à l'écart expliqué et inexpliqué, soit en graphique, soit en tableau.

```
round(results$twofold$variables[[2]][,2:5] ,3)
```

	coef(explained)	se(explained)	coef(unexplained)	se(unexplained)
(Intercept)	0.000	0.000	0.106	0.009
exp_mtra	-0.007	0.001	0.342	0.013
exp_mtra2resc	0.005	0.001	-0.121	0.008
ddipl6	-0.001	0.000	0.000	0.000
ddipl5	0.010	0.000	-0.014	0.001
ddipl4	-0.005	0.000	-0.006	0.001

ddipl3	-0.020	0.001	-0.013	0.001
ddipl1	-0.019	0.001	0.007	0.002
tpartiel	0.168	0.002	-0.050	0.002
secteur0Q	0.027	0.001	-0.028	0.002
secteurBE	0.009	0.000	-0.003	0.000
secteurRU	0.005	0.000	0.006	0.001
secteurFZ	0.003	0.000	0.000	0.000
secteurMN	0.000	0.000	-0.001	0.001
secteurAZ	-0.001	0.000	0.000	0.000
secteurKZ	-0.002	0.000	0.002	0.000
secteurJZ	0.002	0.000	-0.001	0.000
secteurLZ	0.000	0.000	-0.001	0.000
ancentr44	0.001	0.000	-0.083	0.004
ancentr43	0.000	0.000	-0.016	0.001
ancentr42	0.000	0.000	-0.007	0.001

```
plot(results, decomposition = "twofold", group.weight = 1)
```



La partie supérieure du graphique présente la contribution de chaque variable à l'écart expliqué. La variable qui contribue le plus positivement à l'écart expliqué est l'indicatrice de temps partiel, avec une contribution de 0.009 soit 4.8 % de l'écart expliqué total. Autrement dit, presque l'intégralité de la différence de salaires moyens entre hommes et femmes tient au fait que les femmes sont plus souvent en emploi à temps partiel. On notera que l'inclusion de certaines variables peut réduire l'écart inexpliqué : c'est par exemple le cas pour certains niveaux de diplôme. En effet, quand les femmes sont dotées de caractéristiques plus favorables en termes de salaire que les hommes, contrôler de ces caractéristiques réduit la part des écarts qui peut être imputée aux X .

La partie inférieure du graphique ventile l'écart inexpliqué par variable : de même qu'il est possible de détailler les contributions de chaque variable à l'effet de composition, on peut aussi obtenir le détail de l'écart inexpliqué. Cependant, des hypothèses supplémentaires et des précautions particulières sont nécessaires pour analyser et interpréter ces résultats détaillés. On renvoie le lecteur à la section 2.3 pour plus d'éléments sur la décomposition détaillée de l'écart inexpliqué.

2 La validité de l'interprétation

2.1 Effet causal d'appartenance à un groupe et discrimination

Les méthodes de décomposition sont fréquemment utilisées dans le but de mesurer une discrimination entre deux groupes, soit une différence de traitement qui n'est dûe qu'au fait d'appartenir à un groupe plutôt qu'à l'autre. Dans ce cas, l'objectif est d'isoler un effet causal d'appartenance au groupe. Sous quelle condition un écart inexplicé peut-il être interprété comme un effet causal de l'appartenance à un groupe plutôt qu'à l'autre - et donc comme une discrimination ?

Encadré 2 : Décompositions, modèle de Rubin, discrimination

Un individu i est doté des caractéristiques X_i . Soit un "traitement" binaire $T : T_i = 0$ si $i \in \mathcal{P}_0$, $T_i = 1$ si $i \in \mathcal{P}_1$. Les *outcomes* (par exemple les salaires) potentiels s'écrivent :

- $Y_i(0)$ pour l'individu i si $T_i = 0$,
- $Y_i(1)$ si $T_i = 1$.

Or, on observe seulement la réalisation de la variable d'intérêt, soit :

$$Y_i = (1 - T_i)Y_i(0) + T_iY_i(1).$$

Si le modèle est linéaire de la forme $\mathbb{E}(Y | X) = X\beta$ et que l'hypothèse $Y_i(0), Y_i(1) \perp T_i | X_i, \forall i$ (indépendance conditionnelle) est vérifiée, alors $\bar{X}_1\hat{\beta}_0$ est un estimateur convergent de $\mathbb{E}(Y(0) | T = 1)$. Alors, la décomposition de Oaxaca-Blinder :

$$\bar{Y}_1 - \bar{Y}_0 = (\bar{X}_1\hat{\beta}_1 - \bar{X}_1\hat{\beta}_0) + (\bar{X}_1 - \bar{X}_0)\hat{\beta}_0$$

peut être vue comme la contrepartie empirique de :

$$\mathbb{E}(Y(1) | T = 1) - \mathbb{E}(Y(0) | T = 1) + \mathbb{E}(Y(0) | T = 1) - \mathbb{E}(Y(0) | T = 0).$$

La mesure de discrimination (écart inexplicé dans Oaxaca-Blinder) correspond ainsi à l'*average treatment effect on the treated*, soit l'effet du traitement une fois que l'on a contrôlé des différences de caractéristiques entre groupe traité et groupe de contrôle. Cette mesure de discrimination quantifie un effet causal sous l'hypothèse de *conditional independence assumption*.

Pour assimiler l'écart inexplicé à un effet causal, il faut être en mesure d'affirmer qu'aucune différence de caractéristiques inobservées ne subsiste entre les deux groupes, une fois qu'on a contrôlé des caractéristiques observables (encadré 2). C'est une hypothèse forte. Prenons l'exemple des écarts de salaire entre hommes et femmes, lorsque l'on dispose comme variables de contrôle de l'âge, du diplôme et du fait d'être cadre. Une partie de l'écart de salaire entre hommes et femmes est liée aux différences d'âge, de diplôme et de statut entre les hommes et les femmes présents sur le marché du travail. On ne pourra interpréter le reste de l'écart comme de la discrimination que si, pour chaque niveau d'âge, de diplôme et de statut, les hommes et les femmes ont bien un niveau de compétences, y compris inobservées, identique. C'est l'hypothèse

d'indépendance conditionnelle, qui sera formalisée plus bas (section 1.1).

Plusieurs raisons peuvent conduire à ce qu'elle ne soit pas vérifiée. Premièrement, s'il existe une variable omise, qui ne prend pas les mêmes valeurs dans un groupe ou dans l'autre à caractéristiques observables données. L'expérience effective sur le marché du travail pourrait par exemple être plus élevée, à âge donné, chez les hommes que chez les femmes. Dans ce cas, l'écart inexplicé sur-estime le niveau de discrimination car il est en réalité gonflé par une composante qui devrait appartenir à l'écart explicé. Deuxièmement, en présence d'une sélection différenciée sur le marché du travail : si les femmes accèdent plus difficilement à l'emploi que les hommes, les femmes sélectionnées sur le marché du travail pourraient avoir une motivation plus forte que les hommes d'âge et diplôme identiques, motivation qui ne serait pas rétribuée, ou dont la rétribution serait à tort attribuer à d'autres caractéristiques. Dans un tel cas, l'écart attribuable à de la discrimination sera sous-estimé. Enfin, en cas de sélection différenciée dans la CS : si les femmes sont plus rigoureusement sélectionnées pour accéder au statut de cadre, et qu'on contrôle par le fait d'être cadre, on pourra conclure à l'absence de discrimination alors même que les femmes ont une motivation plus grande à niveau d'observables données.

Ces limites de la validité de l'hypothèse d'indépendance conditionnelle doivent être prises en compte dans le choix des variables explicatives. Il y a ainsi un équilibre à trouver en pratique entre l'introduction de contrôles ayant un pouvoir explicatif important et/ou qui sont intéressants pour l'analyse, et la prudence quant aux facteurs qui pourraient fragiliser la condition d'identification. Il faut donc être attentif à ne pas "trop" contrôler et à questionner le choix des variables explicatives incluses dans le modèle : est-ce que pour l'ensemble des X introduits la comparaison des deux groupes a bien un sens ? En général, les variables résultant d'un choix de l'individu doivent être utilisées avec précaution. Un procédé utile lorsqu'on a recours de telles variables est d'introduire les explicatives au fur et à mesure : on commence par les *pre-market factors* - les caractéristiques des individus déterminées avant leur entrée sur le marché du travail-, puis on ajoute les variables de choix comme la CS. On peut ainsi présenter les deux décompositions et préciser que dans la deuxième il est difficile d'assimiler l'écart inexplicé à une discrimination.

Par ailleurs les cas suivants, peu ou pas pertinents dans le cas hommes/femmes, peuvent être rencontrés et rendre invalide l'hypothèse d'indépendance conditionnelle :

- Le fait que l'appartenance au groupe soit le résultat d'une décision de l'individu, par exemple si l'on cherche à étudier les écarts entre public et privé ou encore entre groupes définis selon leur lieu de résidence. Ainsi, les salariés qui choisissent de travailler dans le secteur privé y ont un intérêt plus grand (une espérance de salaires plus élevées par exemple), ce qui se traduit par des inobservables différents. De même, les individus résidant près des zones d'emploi pourraient être plus motivés à niveau de caractéristiques observables donné.
- L'inclusion de variables ne mesurant pas le même phénomène selon le groupe considéré : par exemple lorsque l'on compare immigrés et non-immigrés, ou deux pays dans le cadre d'une comparaison internationale, la variable de diplôme ne reflète pas nécessairement le même niveau de compétences selon le pays dans lequel l'individu a fait ses études.

L'hypothèse d'indépendance conditionnelle autorise que l'effet d'une variable sur le salaire soit mesurée avec biais sur chaque sous-groupe – par exemple l'effet du diplôme sur le salaire capte également l'effet d'une motivation croissante – tant que la structure de corrélation entre diplôme et motivation est la même chez les hommes et chez les femmes (à niveau de diplôme donné, hommes et femmes ont la même motivation)⁴. Attention, cela n'est plus vrai dès lors que l'on cherche à isoler la contribution de chaque variable dans la décomposition détaillée, par exemple connaître la part effectivement liée aux écarts d'éducation dans les écarts de salaire (sans capter par la même occasion la part liée aux écarts de motivation). Cette question sera à nouveau abordée dans la section 2.3.

2.2 Le choix du contrefactuel

Le choix du contrefactuel est crucial, notamment pour bien interpréter les résultats de la décomposition. Dans le cas de l'analyse des inégalités entre une majorité et une minorité, un contrefactuel assez naturel consiste à retenir les caractéristiques du groupe minoritaire et d'y appliquer la structure de salaire du groupe majoritaire. Cela revient implicitement à considérer qu'en l'absence de discrimination salariale entre les deux groupes, tous les salariés seraient rémunérés à la façon dont l'est le groupe en majorité. Les résultats obtenus permettent de répondre à la question de l'existence et de l'ampleur d'une discrimination négative. A l'inverse, en considérant les caractéristiques du groupe majoritaire et en y appliquant les coefficients estimés dans la minorité, on interroge plutôt l'existence de discrimination positive. Enfin, une autre option consiste à raisonner en référence à une moyenne pondérée de $\hat{\beta}_A$ et $\hat{\beta}_B$, ou bien à des coefficients estimés sur l'ensemble de la population avec inclusion d'une indicatrice d'appartenance à l'un des groupes. En procédant ainsi, on tient donc compte de possibles effets d'équilibre. Cela peut par exemple être pertinent pour étudier des inégalités de genre : en l'absence de discrimination, les femmes ne seraient sans doute pas payées de la même manière que le sont les hommes sur un marché du travail avec discrimination de genre.

Encadré 3 : Questions de support commun

- **Dans le cas d'une variable continue :** Si certaines valeurs ne sont pas prises par l'un des groupes, la régression linéaire conduit à "extrapoler" pour les valeurs hors support commun
- **Dans le cas d'une variable catégorielle :** Il faut que chacune des modalités soit connue par chacun des groupes
 - à nuancer si on n'a pas besoin de la décomposition détaillée pour Δ_S , auquel cas on a besoin uniquement d'estimer les $\hat{\beta}_B \rightarrow$ il faut que les B prennent chacune des modalités,

4. Ainsi on autorise au total des différences de caractéristiques inobservées (de motivation par exemple) entre les deux groupes, tant que ces différences sont uniquement liées aux différences de caractéristiques observables (les plus diplômés sont plus motivés, or l'un des groupes est plus diplômé).

— *eg. hommes exerçant le métier de maïeuticien* → d'autant plus problématique que la taille de l'échantillon est réduite.

Variabiles définies dans les deux groupes : une variable non définie pour l'un des groupes ne peut pas être utilisée dans une décomposition. *Si on considère immigrés vs. natifs français, il est problématique d'introduire l'année d'arrivée en France*

Variabilité dans chacun des groupes : une variable constante pour l'un des groupes ne pourra pas non plus être utilisées dans une décomposition. *Si on considère immigrés vs. natifs français, il est problématique d'introduire le pays de naissance.*

2.3 La validité de la décomposition détaillée

On a évoqué précédemment la possibilité, comme pour l'effet de composition, de détailler terme à terme les contributions de chaque variable à l'écart inexpliqué $\hat{\Delta}_S$:

$$\hat{\Delta}_S = \sum_{k=0}^K \hat{\Delta}_{S_k},$$

où pour chaque variable explicative X_k dont la constante, $\hat{\Delta}_{S_k}$ correspond à sa contribution à l'écart inexpliqué, autrement dit :

$$\hat{\Delta}_{S_k}^\nu = \overline{X_{Ak}} \left(\hat{\beta}_{Bk} - \hat{\beta}_{Ak} \right).$$

Néanmoins, cette extension des méthodes de décomposition n'est valide que sous certaines hypothèses et avec des réserves quant à son interprétation.

2.3.1 Une hypothèse plus forte pour l'identification de la décomposition détaillée

La décomposition agrégée peut-être réalisée sans hypothèse sur la forme fonctionnelle du modèle, tant que la distribution conditionnelle des erreurs est la même dans les deux groupes étudiés. En revanche, pour procéder à une décomposition détaillée, il est nécessaire de formuler des hypothèses supplémentaires, afin d'identifier le rôle des $(X_k)_{k=1..K}$ à la fois dans $\hat{\Delta}_S$ et $\hat{\Delta}_X$. Si l'on veut pouvoir attribuer une part de l'écart à une covariable X_k précisément, on revient à l'hypothèse classique sous-jacente à l'estimation sans biais des β dans les équations linéaires initiales : l'hypothèse d'espérance conditionnelle nulle.

2.3.2 Le problème de la modalité omise dans la décomposition détaillée de l'écart inexpliqué

Lorsque certaines caractéristiques X sont catégorielles, la décomposition détaillée de l'écart inexpliqué peut être difficile à interpréter. En effet, les composantes de la part inexpliquée peuvent varier suivant la catégorie de référence omise dans la régression : pour une variable X_k , les parts de Δ_S^ν attribuées à β_0 et à β_k varient. Cette difficulté peut aussi apparaître pour une variable continue dont le zéro n'aurait pas d'interprétation naturelle. Il n'existe pas de solution générale

au problème : un arbitrage entre interprétabilité et comparabilité doit être tranché.

Ainsi, la décomposition détaillée de l'écart inexpliqué peut être satisfaisante si la comparaison au groupe omis a un sens économique (par exemple, les moins qualifiés sont retenus comme référence de la variable diplôme). Dans le cas contraire, l'exercice peut se révéler infructueux. Pour le voir, prenons le cas où le salaire est fonction seulement d'une constante et du secteur (1= services, 0= industrie) :

$$Y_{i,g} = a_g + b_g \text{SECT}_i + \epsilon_i$$

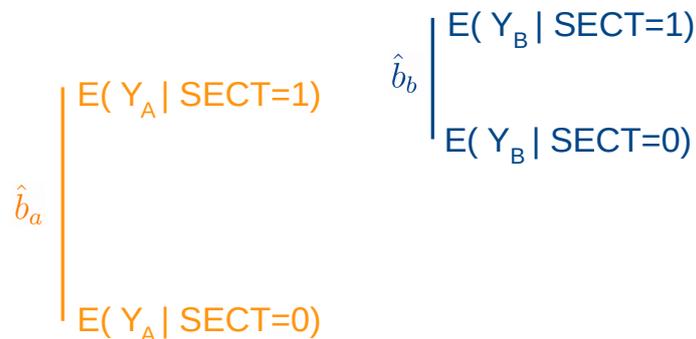
où les estimateurs \hat{a}_g et \hat{b}_g vérifient les relations : $\hat{a}_g = \mathbb{E}(Y_g | \text{SECT} = 0)$ et $\hat{a}_g + \hat{b}_g = \mathbb{E}(Y_g | \text{SECT} = 1)$. L'écart inexpliqué peut se décomposer entre un écart de situation entre hommes et femmes dans l'industrie, soit un terme

$$\hat{\Delta}_{S(\text{constante})} = \hat{a}_b - \hat{a}_a = \mathbb{E}(Y_b - Y_a | \text{SECT} = 0),$$

et un écart d'écart de situation entre hommes et femmes et entre secteurs d'activité,

$$\hat{\Delta}_{S(\text{secteur})} = (\hat{b}_b - \hat{b}_a) \overline{\text{SECT}}_a = [\mathbb{E}(Y_b - Y_a | \text{SECT} = 1) - \mathbb{E}(Y_b - Y_a | \text{SECT} = 0)] \overline{\text{SECT}}_f.$$

FIGURE 2 – Cas de figure où le signe de la contribution à l'inexpliqué est contre-intuitif



Comme l'illustre la figure 2, on peut avoir une situation où les femmes sont désavantagées à la fois dans les services et dans l'industrie, mais où l'écart étant plus grand entre secteurs chez les hommes ($\hat{b}_a > \hat{b}_b$), la contribution de l'appartenance au secteur de l'industrie à l'écart inexpliqué est négatif ($\Delta_{S(\text{SECT})}^\nu < 0$).

Les résultats de la décomposition détaillée de l'inexpliqué sont modifiés par un changement de modalité de référence, avec un transfert entre la contribution à l'inexpliqué de la constante et celle de la variable X . Pour Jones and Kelley (1984), l'interprétation détaillée de la décomposition n'a de sens que pour des variables catégorielles ayant une modalité de référence naturelle. Cette question d'identification a aussi été discutée par Oaxaca and Ransom (1999); Gardeazabal and Ugidos (2004) ou encore Yun (2005, 2008), qui proposent de procéder à une normalisation des

coefficients pour éliminer de la constante l'effet de la modalité omise, par exemple en contraignant à zéro la somme des coefficients de la variable catégorielle. Yun (2005) propose lui de considérer la moyenne des contributions obtenues pour chaque modalité de référence possible associée à chaque variable catégorielle du modèle. La solution proposée par Yun peut être utilisée simplement dans le package `Oaxaca` en renseignant les modalités (sauf une) d'une variable catégorielles de la façon suivante, dans l'appel de la fonction `oaxaca` :

```
results <- oaxaca(formula = logsal ~ exp_mtra + exp_mtra2resc
  + ddip16 + ddip15 + ddip14 + ddip13 + ddip11
  + tpartiel + secteurOQ + secteurBE + secteurRU
  + secteurFZ + secteurMN + secteurAZ + secteurKZ
  + secteurJZ + secteurLZ + ancentr44 + ancentr43
  + ancentr42 | sex | ddip16 + ddip15 + ddip14
  + ddip13 + ddip11, data = data, R=2)
```

Malgré cette correction, la difficulté de l'interprétation n'est pas levée, ce qui rend l'utilisation de la décomposition détaillée de la part inexplicé délicate.

Encadré 4 : Application à la décomposition de différences d'effets fixes

Les méthodes de décomposition peuvent aussi être utilisées dans des modèles à effets fixes. C'est particulièrement utile en économie du travail où l'on s'attend à ce que les caractéristiques inobservables, à la fois des salariés et des entreprises, jouent un rôle essentiel dans les inégalités (Abowd et al., 1999; Lentz and Mortensen, 2010).

Card et al. (2016) s'inspirent des décompositions à la Oaxaca-Blinder pour séparer ce qui, dans l'influence des entreprises sur les inégalités salariales hommes-femmes, provient de la ségrégation des hommes et des femmes dans certaines entreprises et ce qui provient du fait qu'une même entreprise ne rémunère pas de la même manière ses salariés hommes et ses salariées femmes, même si leurs caractéristiques individuelles (compétences) sont identiques. Pour cela, ils proposent un modèle à doubles effets fixes dans l'esprit d'Abowd et al. (1999), pour lequel deux effets fixes sont associés à chaque entreprise, le premier représentant la "prime" que cette entreprise verse à ses salariés hommes et le deuxième celle qu'elle verse à ses salariés femmes. Ces "primes" définissent comment le partage de la richesse se fait au sein de chaque entreprise indépendamment des compétences individuelles des salariés.

Soit le (log) salaire d'un individu à la date t , de sexe $G(i) = g \in \{F, M\}$ et travaillant à la date t dans l'entreprise $J(i, t)$:

$$w_{it}^{G(i)} = \alpha_i + X'_{it}\beta^{G(i)} + \psi_{J(i,t)}^{G(i)} + r_{it}, \quad (5)$$

avec r_{it} composé d'un terme d'erreur individuel et des éléments variant dans le temps du surplus de l'entreprise.

Une telle écriture décompose donc le salaire en fonction d'un effet fixe individuel α_i , d'un effet entreprise pour les hommes et pour les femmes $\psi_{J(i,t)}^{G(i)}$, et de covariables aux rendements spécifiques pour les hommes et pour les femmes. Avec $\psi_{J(i,t)}^g$ l'effet fixe spécifique pour l'entreprise $J(i,t)$ pour le genre g , on peut réécrire l'écart moyen entre effets entreprises moyens des hommes et des femmes de la manière suivante :

$$\begin{aligned} \mathbb{E} \left[\psi_{J(i,t)}^M \mid g = M \right] - \mathbb{E} \left[\psi_{J(i,t)}^F \mid g = F \right] &= \underbrace{\mathbb{E} \left[\psi_{J(i,t)}^M - \psi_{J(i,t)}^F \mid g = M \right]}_{\text{Effet bargaining}} \\ &+ \underbrace{\mathbb{E} \left[\psi_{J(i,t)}^F \mid g = M \right] - \mathbb{E} \left[\psi_{J(i,t)}^F \mid g = F \right]}_{\text{Effet sorting}} \end{aligned}$$

Le premier terme de cette décomposition correspond à la différence de l'effet fixe entreprise moyen chez les hommes et chez les femmes, *si les femmes travaillaient dans les mêmes entreprises que les hommes*, soit la différence, pour un même effet fixe entreprise moyen, de captation du surplus de l'entreprise par les hommes et par les femmes (ou *bargaining effect*). Le second élément de la décomposition correspond à la différence entre l'effet entreprise moyen pour les femmes *si elles travaillaient dans les mêmes entreprises que les hommes*, et leur véritable effet entreprise moyen, étant donné leur répartition dans les entreprises, soit à la pénalité salariale liée au fait que les femmes travaillent dans des entreprises qui paient moins bien leurs salariés, toutes choses égales par ailleurs (ou effet de *sorting*). La méthode de décomposition est ici utilisée comme un outil d'identification des composantes d'inégalités intra-entreprises et inter-entreprises entre hommes et femmes.

3 Variable d'intérêt dichotomique et écart entre proportions

On a jusqu'ici étudié l'écart entre deux groupes selon une variable continue, par exemple le salaire. De nombreuses variables sont toutefois dichotomiques : si l'on considère par exemple le fait d'être au chômage, on sera amené à décomposer l'écart de taux de chômage entre les deux sous-populations considérées.

3.1 La décomposition d'Oaxaca-Blinder pour une variable dichotomique

La décomposition présentée jusqu'ici pour le cas d'une variable Y continue peut être directement transposée à une variable d'intérêt dichotomique, dès lors qu'il est raisonnable de modéliser celle-ci par une régression linéaire. Cela présente l'avantage de la simplicité : par exemple, lorsqu'on souhaite comparer l'effet qu'exerce une variable explicative sur la probabilité de chômage dans un groupe et dans l'autre, on peut simplement estimer les modèles linéaires dans les deux groupes puis comparer les coefficients sans qu'il soit nécessaire de recalculer des effets marginaux. Par ailleurs, la décomposition détaillée variable par variable est immédiate (voir section 2.3) ce qui n'est plus le cas dès lors qu'on s'éloigne du modèle linéaire.

L'approximation linéaire peut toutefois s'avérer problématique lorsque les différences de caractéristiques observables entre les deux groupes sont marquées, ou que l'événement modélisé est très rare ou au contraire très courant, car on court alors le risque de s'appuyer implicitement sur un contrefactuel dénué de sens (par exemple un taux de chômage négatif). On peut alors se tourner vers une modélisation non-linéaire de type logit ou probit.

3.2 Modèle de Fairlie

Fairlie (2005) a adapté la décomposition d'Oaxaca-Blinder au cas d'une variable dichotomique en recourant à un modèle probit ou logit. Les trois étapes de la méthode de Fairlie sont les suivantes, en prenant le groupe B comme référence et comme variable d'intérêt le fait d'être au chômage (qui vaut 0 ou 1) :

- (1) On modélise la probabilité d'être au chômage au sein du groupe B de la façon suivante : $P(Y_B = 1|X) = F(X\beta_B)$, où $F(\cdot)$ est la fonction de répartition de la loi normale (modèle probit) ou de la loi logistique (modèle logit).
- (2) On calcule alors, pour chaque individu i du groupe A , sa probabilité prédite d'être au chômage en appliquant $F(X\beta_B)$ aux caractéristiques observables de i . C'est-à-dire que pour chaque femme, on calcule sa probabilité d'être au chômage si ses caractéristiques restaient inchangées mais étaient valorisées (les exposaient au chômage) comme celles des hommes. Pour une femme avec des caractéristiques X_i , on obtient donc $\hat{P}_B(Y_i = 1|X_i) = F(\hat{\beta}_{B0} + \sum_{k=1}^K X_{ik}\hat{\beta}_{Bk})$.
- (3) On effectue la moyenne de ces probabilités prédites pour l'ensemble des individus du groupe A : $\frac{1}{N_A} \sum_{i \in A} \hat{P}_B(Y_i = 1|X_i)$.

Comme pour la décomposition d'Oaxaca-Blinder, on obtient ainsi un contrefactuel répondant à la question suivante : quel serait le taux de chômage des individus du groupe A si leurs caractéristiques étaient valorisées de la même manière que pour le groupe B ? Ces étapes permettent d'obtenir la décomposition agrégée, qui s'écrit :

$$\underbrace{\frac{1}{N_B} \sum_{i \in B} \mathbf{1}_{\{Y_i=1\}} - \frac{1}{N_A} \sum_{i \in A} \hat{P}_B(Y_i = 1|X_i)}_{\text{Effet de composition (lié aux X)}} + \underbrace{\frac{1}{N_A} \sum_{i \in A} \hat{P}_B(Y_i = 1|X_i) - \frac{1}{N_A} \sum_{i \in A} \mathbf{1}_{\{Y_i=1\}}}_{\text{Ecart inexplicé (à X donnés)}}$$

Encadré 5 : La décomposition détaillée au-delà du cas linéaire

Deux propriétés sont particulièrement souhaitables pour la décomposition détaillée : l'additivité et l'invariance à l'ordre. On entend par additivité le fait que les contributions à l'expliqué de chaque variable se somment bien en la part expliquée totale, autrement dit : $\Delta_X^\nu = \sum_{k=1}^K \Delta_{X_k}^\nu$. Cette propriété est satisfaite dans le cadre linéaire simple mais elle n'est pas forcément garantie hors de celui-ci, par exemple dans les approches présentées dans la section 4.

Elle sera généralement satisfaite dans une procédure séquentielle consistant à remplacer la

distribution de X_1 puis de X_2 etc., jusqu'à ce que la distribution des X ait été entièrement remplacée. Mais comme l'impact du changement d'une variable donnée dépend généralement de la distribution des autres variables, on peut alors avoir une décomposition détaillée qui dépend de l'ordre dans lequel on la réalise. L'invariance à l'ordre n'est donc pas respectée.

Si la décomposition agrégée du modèle de Fairlie est facile à obtenir, la version détaillée est nettement plus difficile à calculer. En effet, dans la construction du contrefactuel pour la décomposition agrégée $F(\hat{\beta}_{B0} + \sum_{k=1}^K X_{ik}\hat{\beta}_{Bk})$ implicitement on "remplace" les X des hommes par les X des femmes. Pour détailler l'effet de composition il faut donc remplacer successivement chaque X_k des hommes par les X_k des femmes. Prenons l'exemple d'un cas à trois variables explicatives. Pour calculer la contribution de X_3 à l'effet de composition, il faut prendre la différence entre :

$$F(\hat{\beta}_{0B} + X_{1B}\hat{\beta}_{1B} + X_{2B}\hat{\beta}_{2B} + X_{3B}\hat{\beta}_{3B})$$

et $F(\hat{\beta}_{0B} + X_{1B}\hat{\beta}_{1B} + X_{2B}\hat{\beta}_{2B} + \underline{X_{3A}}\hat{\beta}_{3B}).$

Pour un homme à X_1 et X_2 donnés, par quelle valeur remplace-t-on son X_3 ? On voudrait une sorte d'appariement entre les hommes et les femmes pour tenir compte de la structure de corrélation entre les variables X_1 , X_2 et X_3 . Pour ce faire, la solution proposée par Fairlie (2005) consiste en quatre étapes, par exemple si la variable d'intérêt est la probabilité d'être au chômage pour les actifs :

- (1) On tire un échantillon dans la population majoritaire, de même taille que celle de la population minoritaire.
- (2) Au sein de chaque échantillon, on classe les individus selon leur propension à être au chômage.
- (3) On apparie l'homme ayant la plus forte propension à être au chômage à la femme ayant la plus forte propension à être au chômage, etc.
- (4) Pour un homme donné, on remplace la valeur de X_k considérée par celle prise par l'individu femme apparié.

On reproduit ces étapes un grand nombre de fois, en tirant à chaque fois un nouvel échantillon.

Cette procédure est très intensive en calcul, et elle ne résout le problème de l'impossibilité d'avoir une décomposition détaillée additive et non sensible à l'ordre. Pour une approche plus simple, on pourra préférer l'approximation de Yun (2004) : celle-ci consiste à repartir de l'effet de composition agrégé estimé selon Fairlie et à le désagréger selon un système de poids attribuant à chaque variable le poids $\frac{(\bar{X}_B^k - \bar{X}_A^k)\hat{\beta}_B^k}{\sum_k (\bar{X}_B^k - \bar{X}_A^k)\hat{\beta}_B^k}$ avec les $\hat{\beta}_B^k$ estimés par logit ou probit. Cette méthode peut cependant poser problème lorsque les prédictions se prêtent mal à l'approximation linéaire, typiquement quand elles sont hors de l'intervalle entre 0 et 1 et/ou lorsqu'il existe de fortes différences dans les X entre les deux groupes.

3.3 Décomposition à la Fairlie dans R

On propose d'appliquer la méthode de Fairlie aux écarts de probabilité d'accès à des postes d'encadrement entre hommes et femmes, toujours à partir des données de l'enquête emploi. On cherche à expliquer ces différences par des différences de caractéristiques : ancienneté dans l'entreprise, quotité de temps de travail (entre 1 et 6, ou 6=temps plein et 1=moins de 50 %), expérience potentielle et diplôme.

```
#on garde les actifs pour lesquels on dispose du salaire
base_empl<-data[ data$acteu==1 & !is.na(data$sal),]

# Quelle proportion d'hommes occupe des fonctions d'encadrement ?
prop.h<-mean(base_empl$encadr[base_empl$sexe=="1"])
prop.h

## [1] 0.251

# Et parmi les femmes ?
prop.f<-mean(base_empl$encadr[base_empl$sexe=="2"])
prop.f

## [1] 0.136

# Ecart hommes-femmes
prop.h-prop.f

## [1] 0.115
```

On estime pour chaque sexe un modèle logistique d'accès à des fonctions d'encadrement :

```
logitH<-glm(encadr ~ as.factor(ancentr4) + as.factor(quotite)
+ exp+as.factor(ddipl), family=binomial (link='logit'),
data=base_empl[base_empl$sexe=="1",])

logitF<-glm(encadr ~ as.factor(ancentr4) + as.factor(quotite)
+ exp+as.factor(ddipl), family=binomial (link='logit'),
data=base_empl[base_empl$sexe=="2",])
```

Mettons que l'on souhaite connaître la probabilité contrefactuelle d'encadrement parmi les hommes, s'ils avaient les caractéristiques des femmes (ou dit autrement, la probabilité contrefactuelle d'encadrement des femmes si leurs caractéristiques étaient valorisées comme celles des hommes). Il suffit pour cela de calculer les prédictions individuelles selon le modèle des hommes, puis de calculer la moyenne de ces probabilités prédites parmi les femmes.

```

base_empl$pH<-predict(logitH,base_empl,type='response')
prop.cf<-mean(base_empl$pH[base_empl$sexe=="2"],na.rm=TRUE)
prop.cf

## [1] 0.231

#Ecart expliqué
expl<-prop.h - prop.cf
expl

## [1] 0.0199

#Ecart inexpliqué
inexpl<-prop.cf - prop.f
inexpl

## [1] 0.0951

```

Ce contrefactuel nous permet de mesurer un écart expliqué et un écart inexpliqué : l'écart expliqué se calcule comme l'écart entre la probabilité contrefactuelle et la proportion d'encadrement mesurée parmi les hommes, car celui-ci provient bien uniquement de différences de caractéristiques. Au contraire pour l'écart inexpliqué, on raisonne à caractéristiques données (celles des femmes). On peut par ailleurs s'assurer que le modèle permet bien de reconstituer la probabilité effective d'encadrement des hommes, à travers le calcul de $\frac{1}{N_B} \sum_{i \in B} F(\hat{\beta}_B X_i)$.

```

mean(base_empl$pH[base_empl$sexe=="1"], na.rm=TRUE)

## [1] 0.251

```

On peut effectuer la même décomposition en repartant, non pas d'un logit, mais d'un modèle de probabilité linéaire (ie. simples MCO). A noter : ce code permet d'obtenir plus généralement la décomposition d'Oaxaca-Blinder agrégée, quelle que soit la variable considérée.

```

lpmH<-lm(encadr ~ as.factor(ancentr4) + as.factor(quotite)
          + exp+as.factor(ddipl), data=base_empl[base_empl$sexe=="1",])
base_empl$pH<-predict(lpmH,base_empl)
mean(base_empl$pH[base_empl$sexe=="2"], na.rm=TRUE)

## [1] 0.237

```

On voit toutefois que dans d'assez nombreux cas, la probabilité prédite est en dehors de [0,1].

```

table(base_empl$pH<0 | base_empl$pH>1)

##
## FALSE TRUE
## 191191 7344

```

Pour la décomposition détaillée, on va procéder à l'approximation de Yun (voir encadré 5). Pour cela, on a besoin de récupérer le vecteur des $(\bar{X}_B^k - \bar{X}_A^k)\hat{\beta}_B^k, k = 1...K$, qu'on nomme ci-dessous `delta.X.beta`.

```
#On récupère le vecteur des coefficients chez les hommes,
#ainsi que le vecteur des X moyens dans les deux groupes
coeffs.H<-logitH$coefficients

X.H <- model.matrix(~ as.factor(ancentr4) + as.factor(quotite)
                    + exp+as.factor(ddipl), data=base_empl[base_empl$sexe=="1",])
X.moy.H<-apply(X.H,2,mean)

X.F <- model.matrix(~ as.factor(ancentr4) + as.factor(quotite)
                    + exp+as.factor(ddipl), data=base_empl[base_empl$sexe=="2",])
X.moy.F<-apply(X.F,2,mean)

#On calcule alors delta.X.beta
delta.X.beta<-(X.moy.H- X.moy.F)*coeffs.H
delta.X.beta

##          (Intercept) as.factor(ancentr4)2 as.factor(ancentr4)3 as.factor(ancentr4)4
##          0.00000      -0.00164          -0.00234          -0.00338
## as.factor(ancentr4)5 as.factor(quotite)2 as.factor(quotite)3 as.factor(quotite)4
##          -0.02024      -0.01253          -0.03157          -0.02500
## as.factor(quotite)5 as.factor(quotite)6          exp      as.factor(ddipl)3
##          -0.01000          0.37123          -0.00951          0.01506
## as.factor(ddipl)4   as.factor(ddipl)5   as.factor(ddipl)6   as.factor(ddipl)7
##          0.01051          -0.07938          0.00455          -0.03824

# Part liée à l'ancienneté :
# 4 modalités qui correspondent aux éléments 2 à 5 de delta.X.beta
part.ancentr<-expl*sum(delta.X.beta[2:5])/sum(delta.X.beta)

#Idem pour quotité (6 à 10), expérience potentielle (11) et diplôme (12 à 16)
part.quotite<-expl*sum(delta.X.beta[6:10])/sum(delta.X.beta)
part.exp<-expl*sum(delta.X.beta[11])/sum(delta.X.beta)
part.ddipl<-expl*sum(delta.X.beta[12:16])/sum(delta.X.beta)
```

4 Décompositions au-delà de la moyenne

Lorsque la variable d'intérêt Y est continue, on la résume souvent par sa moyenne : on cherche alors à expliquer l'écart entre moyennes calculées pour chacun des deux groupes. On peut toutefois souhaiter aller "au-delà de la moyenne" et s'intéresser à des inégalités en certains endroits

de la distribution de Y , ou plus généralement à d'autres statistiques que la moyenne : en termes de salaires par exemple, il peut exister un phénomène de type *plafond de verre* lorsqu'un des deux groupes ne parvient pas aux salaires les plus élevés. Dans ce cas, il sera plus pertinent de s'intéresser au sommet de la distribution des salaires, plutôt qu'au salaire moyen. De même, lorsqu'on effectue une comparaison intertemporelle ou internationale, c'est souvent à une statistique caractérisant les inégalités que l'on s'intéresse (par exemple écart interdécile, le coefficient de Gini, etc.), pour chaque période ou pour chaque pays, plutôt qu'à la seule moyenne.

Dans ce cas plus général, on va donc s'intéresser à l'écart entre la distribution de Y observée dans le groupe A , et celle observée dans le groupe B . Pour ce faire, on va employer la notion de *distribution conditionnelle*, qui correspond simplement à la fonction qui associe à un ensemble de caractéristiques X , la distribution que prend Y pour chaque valeur de ces caractéristiques. Si par exemple on considère une unique variable binaire X (le fait d'être cadre ou non), et qu'on s'intéresse à la distribution des salaires dans le groupe A , la distribution conditionnelle de Y à X dans le groupe A – qu'on notera $F_{Y_A|X}$ – associe à $X = 1$ la distribution des salaires parmi les cadres du groupe A , et à $X = 0$ la distribution des salaires parmi les non-cadres du groupe A . De façon générale, en considérant un ensemble de caractéristiques X plus vaste, on peut écrire la distribution des salaires effectivement observée dans le groupe A (la distribution *non-conditionnelle* F_{Y_A} , qu'on pourra également noter $F_{Y_A|X_A}$) comme la résultante de la distribution conditionnelle $F_{Y_A|X}$, appliquée à la répartition des caractéristiques X dans le groupe A . On a ainsi⁵ :

$$F_{Y_A}(= F_{Y_A|X_A}) = \int F_{Y_A|X}(y|x)dF_{X_A}(x)$$

avec F_{X_A} la distribution des caractéristiques observables dans le groupe A . La figure 3 met en avant le passage, pour chacun des groupes A et B , entre distributions des caractéristiques observables X (que celles-ci soient discrètes ou continues), distributions conditionnelles de Y qui valorisent ces caractéristiques, et distributions non-conditionnelles. L'écart entre les distributions observées F_{Y_A} et F_{Y_B} peut ainsi trouver deux sources : un écart entre distribution des caractéristiques observables F_{X_A} et F_{X_B} , ou un écart entre distributions conditionnelles $F_{Y_A|X}$ et $F_{Y_B|X}$, c'est-à-dire entre valorisations des X en termes de distributions de salaires.

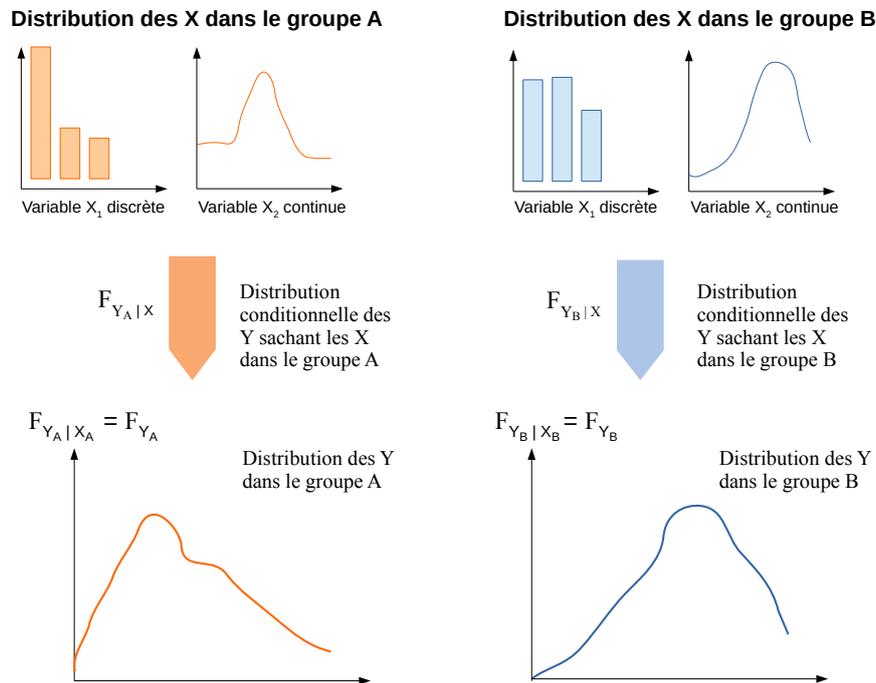
Afin de distinguer entre effet de composition, et écart inexplicé à caractéristiques données, on va introduire un terme correspondant par exemple à la distribution contrefactuelle qui s'appliquerait si les caractéristiques présentes dans le groupe A étaient valorisées comme chez les B :

$$F_{Y_B|X_A} = \int F_{Y_B|X}(y|x)dF_{X_A}(x)$$

Supposons que l'on s'intéresse à une statistique ν de la distribution en particulier, par exemple le dernier décile : on souhaiterait décomposer l'écart entre le dernier décile de salaire dans le groupe B , et le dernier décile de salaire dans le groupe A . On peut décomposer l'écart de ν entre

5. Les notations proposées ici sont légèrement simplifiées par rapport à celles de Fortin et al. (2011). On ne reprend notamment pas l'indicatrice d'appartenance au groupe, $D_g, g = A, B$ et on indice directement les distributions en désignant le groupe concerné.

Figure 3 – Distribution jointe de X et Y dans chaque groupe



groupes B et A de la façon suivante :

$$\nu(F_{Y_B}) - \nu(F_{Y_A}) = \underbrace{[\nu(F_{Y_B|X_B}) - \nu(F_{Y_B|X_A})]}_{\text{Effet de composition}} + \underbrace{[\nu(F_{Y_B|X_A}) - \nu(F_{Y_A|X_A})]}_{\text{Écart inexpliqué}} \quad (6)$$

Le premier terme correspond à l'effet de composition : on voit en effet apparaître un écart lié aux caractéristiques observables (X_A vs. X_B), valorisées dans les deux cas par la même distribution conditionnelle $F_{Y_B|X}$. Pour le deuxième terme au contraire, on raisonne à caractéristiques données (X_A) : il s'agit de l'écart inexpliqué. Plusieurs des méthodes de décomposition de l'écart entre distributions reposent ainsi sur la construction de la distribution contrefactuelle $F_{Y_B|X_A}$.⁶ On peut distinguer deux façons de parvenir à la distribution contrefactuelle $F_{Y_B|X_A}$:

- soit on part de la distribution des salaires dans le groupe B ($F_{Y_B|X_B}$), mais on modifie la distribution de leurs caractéristiques observables de façon à ce qu'elle soit la même que dans le groupe A (on "remplace" ainsi F_{X_B} par F_{X_A}). Cela correspond aux méthodes par repondération (DiNardo, Fortin, and Lemieux, 1996). Ce procédé est représenté dans la partie gauche de la figure 4, et présenté dans la section suivante.
- soit on estime directement la distribution conditionnelle du groupe B ($F_{Y_B|X}$), et on l'applique ensuite aux caractéristiques X du groupe A . Cela correspond aux méthodes d'estimation de la distribution conditionnelle (Chernozhukov et al., 2013; Machado and

6. Comme dans le cas de la décomposition d'Oaxaca-Blinder, d'autres distributions contrefactuelles peuvent bien sûr être envisagées, en premier lieu $\nu(F_{Y_A|X_B})$.

Mata, 2005). Ce procédé est représenté dans la partie droite de la figure 4, et présentée dans l'encadré 6.

4.1 La méthode de repondération

Afin de construire la distribution contrefactuelle $F_{Y_B|X_A}$ correspondant à la distribution des Y du groupe B , si celui-ci présentait les mêmes caractéristiques observables que celles du groupe A , DiNardo et al. (1996) proposent d'ajuster les poids des observations du groupe B afin de rendre leurs caractéristiques observables similaires à celles des individus du groupe A . Par exemple, si l'on souhaite décomposer l'écart entre les distributions de salaire des hommes et des femmes en contrôlant du statut de cadre, et que l'on suppose que les hommes accèdent plus souvent au statut de cadre que les femmes : on va repondérer à la baisse les observations des hommes exerçant des fonctions d'encadrement ; et à la baisse les observations des hommes exerçant des fonctions d'encadrement. Á partir de la distribution des salaires pour les observations hommes ainsi repondérées (qui correspond ici à la distribution contrefactuelle $F_{Y_H|X_F}$), on peut calculer très facilement n'importe quelle statistique ν et parvenir à la décomposition 6. L'étape de calcul des poids de repondération peut elle-même s'effectuer très aisément.

En effet, le facteur de repondération $\Psi(X)$ qui, appliqué à chaque observation du groupe B , permet de rendre la distribution des caractéristiques du groupe B similaire à celle du groupe A s'écrit (en notant $g = A, B$ la variable d'appartenance au groupe) :

$$\Psi_{DFL}(X) = \frac{P(X|g=A)}{P(X|g=B)} = \frac{P(g=A|X) \cdot P(g=B)}{P(g=B|X) \cdot P(g=A)} = \frac{P(g=A|X)}{1 - P(g=A|X)} \cdot \frac{1 - P(g=A)}{P(g=A)}$$

$P(g=A)$ correspond simplement à la proportion d'individus du groupe A dans la population. Afin d'obtenir une estimation de $P(g=A|X)$, on modélise la probabilité d'appartenir au groupe A , sur l'ensemble de l'échantillon, en fonction des caractéristiques observables X . L'estimation peut être faite par logit ou probit⁷. Ce modèle fournit directement pour chaque individu de caractéristiques X , la probabilité prédite d'appartenir au groupe A , c'est-à-dire $\hat{P}(g=A|X)$. On calcule alors le facteur de repondération $\hat{\Psi}_{DFL}(X)$ de façon très simple :

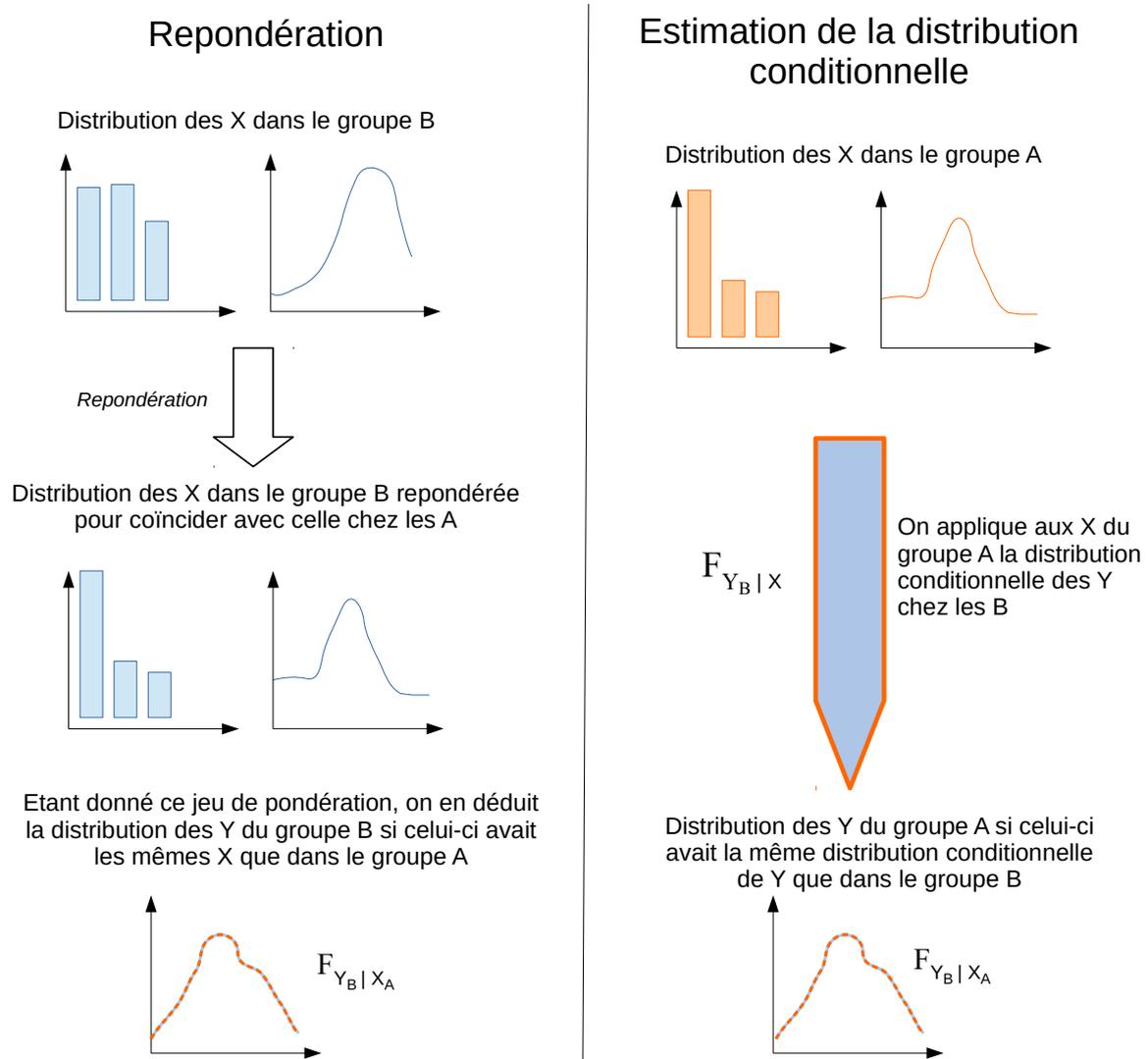
$$\hat{\Psi}_{DFL}(X) = \frac{\hat{P}(g=A|X)}{1 - \hat{P}(g=A|X)} \cdot \frac{1 - \hat{P}(g=A)}{\hat{P}(g=A)}$$

Bien que très simple à mettre en œuvre, cette méthode doit être utilisée avec précaution en cas de problème de support commun, car le facteur de repondération peut alors avoir un comportement erratique. Notamment, si $P(g=B|X) \rightarrow 0$ et $P(g=A|X) \rightarrow 1$, ce qui sera le cas si une caractéristique en particulier est très rare au sein du groupe B relativement au groupe A , $\Psi(X)$ peut devenir très grand pour les individus B détenant cette caractéristique : ces observations repondérées risquent alors de porter à elles seules toute la distribution contrefactuelle⁸.

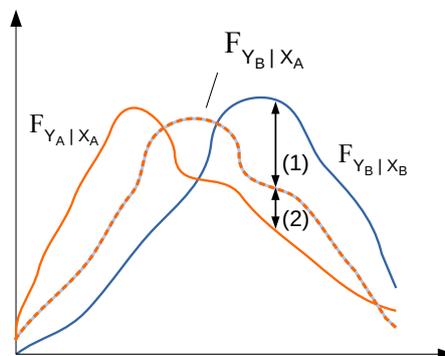
7. Hirano et al. (2003) proposent alternativement l'emploi d'un modèle non-paramétrique, permettant de tenir compte de façon plus flexible de la structure de corrélation entre les variables.

8. Il est ainsi nécessaire de s'assurer, lorsque l'estimation du facteur de repondération pour chacun des B est effectuée, que celui-ci ne prend pas de valeur anormalement élevée ou faible. En pratique, on peut regarder la façon dont les poids après repondération sont distribués.

Figure 4 – Comparaison des méthodes de décomposition au-delà de la moyenne : repondération ou estimation de la distribution conditionnelle



Les deux méthodes de calcul permettent d'estimer la même distribution contrefactuelle, dont on peut ensuite déduire les parts expliquée (1) et inexpliquée (2) de l'écart en tout point de la distribution.



La méthode de repondération initialement proposée par DiNardo et al. (1996) permet d'isoler la participation à l'effet de composition d'une variable binaire. Des travaux ultérieurs, notamment Altonji et al. (2012), ont proposé des extensions à des variables catégorielles ou continues. Toutefois, la décomposition détaillée obtenue est non-additive, si l'on remplace pour chaque X_k la distribution au sein du groupe B par celle du groupe A , tout en conservant pour les autres explicatives la distribution des B . Si l'on procède plutôt de façon séquentielle en remplaçant successivement la distribution de X_1 , puis de X_2 , et ainsi de suite jusqu'à ce que la distribution de l'ensemble des X soit celle du groupe A , la décomposition détaillée obtenue est additive mais dépendante à l'ordre dans lequel on procède.

Application 3 : décomposition par repondération dans R

Considérons ici l'écart entre la distribution des salaires des descendants d'immigrés maghrébins (le groupe A), et celle des non-descendants (le groupe B). On souhaite repondérer les non-descendants pour qu'ils ressemblent, en termes d'expérience potentielle et de diplôme, aux descendants d'immigrés maghrébins. On va ainsi être amené à augmenter le poids des non-descendants dont les caractéristiques sont courantes parmi les descendants (par exemple, les individus jeunes) relativement à celui des non-descendants dont les X sont rares parmi les descendants d'immigrés maghrébins.

La première étape consiste alors à estimer la probabilité conditionnelle à X d'appartenir au groupe des descendants d'immigrés maghrébins, relativement aux non-descendants.

```
logit<-glm(magh ~ exp + exp2 + ddipl, family=binomial (link='logit'),
           data=base)
summary(logit)$coefficients
```

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-3.72760	0.092790	-40.173	0.00e+00
exp	0.05094	0.008162	6.241	4.34e-10
exp2	-0.00247	0.000207	-11.953	6.25e-33
ddipl1	0.02319	0.083515	0.278	7.81e-01
ddipl4	0.37256	0.079151	4.707	2.51e-06
ddipl5	0.35315	0.077245	4.572	4.84e-06
ddipl6	0.73597	0.111853	6.580	4.71e-11
ddipl7	1.08624	0.086576	12.547	4.15e-36

On voit par exemple qu'à expérience donnée, les descendants d'immigrés maghrébins sont $e^{1.09} = 2.96$ fois plus susceptibles d'être sans diplôme plutôt que diplômés d'un Bac+3 ou plus, relativement aux non-descendants. On sera donc amené à repondérer à la hausse les observations des non-descendants sans diplôme.

Le facteur de repondération $\hat{\Psi}(X) = \frac{\hat{P}(g=B|X)}{\hat{P}(g=A|X)} \cdot \frac{\hat{P}(g=A)}{\hat{P}(g=B)}$ est calculé de la façon suivante :

```
p<-predict(logit,type='response')
w1<-ifelse(base$magh==0,
           p/(1-p)*(1-mean(base$magh))/mean(base$magh), 1)
```

On peut s'assurer que cette opération a bien rendu comparable les deux populations selon les dimensions observables considérées. Par exemple, la proportion d'individus sans diplôme était initialement de 0.11 parmi les non-descendants contre 0.158 parmi les descendants. Elle est de 0.158 parmi les non-descendants repondérés.

Une fois les pondérations obtenues, on pourra directement calculer la statistique d'intérêt sur la distribution contrefactuelle, c'est-à-dire ici sur la distribution de salaire des non-descendants repondérés pour ressembler aux descendants en termes de caractéristiques observables. La fonction `wtd.quantile` du *package* `Hmisc` permet notamment de calculer des quantiles en incluant des pondérations. On écrira par exemple, pour obtenir les déciles du log-salaire dans la population des non-descendants repondérés :

```
library(Hmisc)
grid<-seq(0.1,0.9,0.1)
ref<-base$magh==0
dfl.Fc<-wtd.quantile(base$logsal[ref], weights=w1[ref], probs=grid)
```

Les déciles de log-salaire ainsi obtenus pour la distribution contrefactuelle, ainsi que les distributions initiales, sont présentées à la figure 5.

Pour les premiers déciles de revenus, la distribution contrefactuelle est très proche de celle des descendants d'immigrés maghrébins. Les différences d'expérience et de diplôme entre les deux groupes expliquent presque entièrement les différences de salaires mesurées entre ces quantiles. Toutefois, plus on progresse dans la distribution, plus l'écart de salaire inexpliqué entre les deux groupes devient grand.

```
#Ecart total
round(dfl.Fref-dfl.Fmagh,3)

 10%  20%  30%  40%  50%  60%  70%  80%  90%
0.068 0.067 0.069 0.074 0.094 0.107 0.134 0.160 0.188

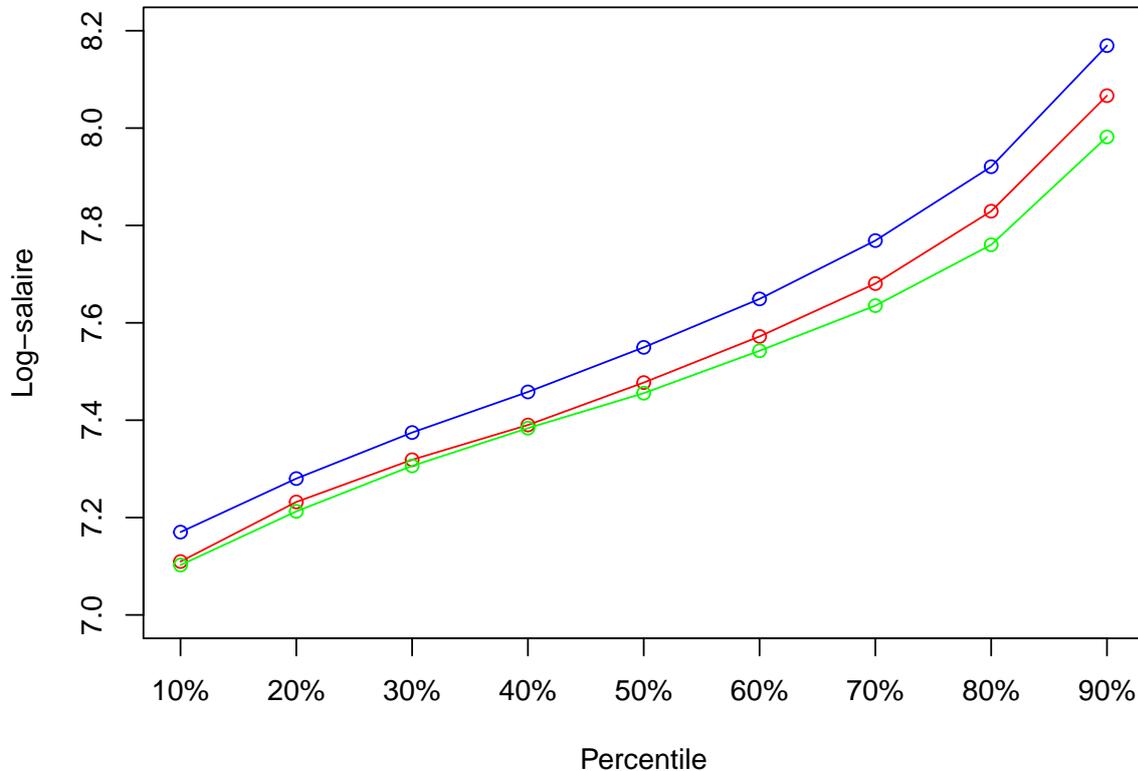
#Dont effet de composition
round(dfl.Fref-dfl.Fc,3)

 10%  20%  30%  40%  50%  60%  70%  80%  90%
0.060 0.048 0.056 0.068 0.072 0.077 0.088 0.091 0.103

#Dont écart inexpliqué
round(dfl.Fc-dfl.Fmagh,3)

 10%  20%  30%  40%  50%  60%  70%  80%  90%
0.007 0.019 0.013 0.007 0.022 0.029 0.045 0.069 0.085
```

FIGURE 5 – Distributions de log-salaire selon le groupe considéré



4.2 Les décompositions par régression quantile non-conditionnelle

Pour détailler le rôle de chacune des variables dans la décomposition de façon à la fois additive et dépendante à l'ordre, Firpo et al. (2007) proposent une solution qui se rapproche de l'esprit de la décomposition d'Oaxaca-Blinder, mais adaptée au cas où l'on considère d'autres statistiques que la moyenne, notamment les quantiles de la distribution. Notons que dans le cas de la moyenne, le modèle de régression linéaire permettra d'écrire la moyenne empirique de Y comme $\bar{Y} = \bar{X}\hat{\beta}$, ce qui autorise ensuite à procéder à la décomposition d'Oaxaca-Blinder. Or si l'on considère le quantile d'ordre τ de la distribution de Y (on le note $Q^\tau(Y)$), il existe une modélisation qui permettra *in fine* d'exprimer le quantile empirique $\hat{Q}^\tau(Y)$ comme une fonction linéaire des X moyens, c'est-à-dire comme $\hat{Q}^\tau(Y) = \bar{X}\hat{\gamma}^\tau$: c'est la méthode des régressions quantiles non-conditionnelles, proposée par Firpo et al. (2009) (cf. Encadré 5).

Les valorisations γ_A^τ et γ_B^τ qu'on estime au sein de chaque sous-population pour un quantile d'ordre τ donné sont l'équivalent des valorisations β_A, β_B dans le cas de la décomposition de la moyenne. On effectuera une décomposition pour chaque $\tau \in [0, 1]$ auquel on s'intéresse (par exemple $\tau = 0.9$ si l'on souhaite se pencher sur le haut de la distribution) – en règle générale, on se penchera sur des points tout au long de la distribution). γ^τ correspond à la valorisation des X en un point τ donné de la distribution, mais c'est bien à la moyenne de X dans toute

la (sous-)population qu'on les applique – même si, par exemple, on estime l'effet d'être "Sans diplôme" au quantile d'ordre $\tau = 0.9$ de la distribution de salaire.

Pour un τ donné, une fois les valorisations γ_A^τ et γ_B^τ estimées, on a d'une part $\widehat{Q}_A^\tau(Y) = \bar{X}_A \hat{\gamma}_A^\tau$, d'autre part $\widehat{Q}_B^\tau(Y) = \bar{X}_B \hat{\gamma}_B^\tau$. Encore une fois, on voit que l'écart entre les quantiles d'ordre τ dans les groupes A et B peut provenir soit d'une différence de caractéristiques X entre les deux sous-populations, soit d'une différence dans la valorisation de ces caractéristiques moyennes en un point donné de la distribution. La décomposition de l'écart entre $\widehat{Q}_B^\tau(Y)$ et $\widehat{Q}_A^\tau(Y)$ s'écrit alors comme :

$$\widehat{Q}_B^\tau(Y) - \widehat{Q}_A^\tau(Y) = \underbrace{\sum_{k=1}^K (\bar{X}_{Bk} - \bar{X}_{Ak}) \hat{\gamma}_{Bk}^\tau}_{\hat{\Delta}_X^\tau} + \underbrace{\hat{\gamma}_{B0}^\tau - \hat{\gamma}_{A0}^\tau + \sum_{k=1}^K \bar{X}_{Ak} (\hat{\gamma}_{Bk}^\tau - \hat{\gamma}_{Ak}^\tau)}_{\hat{\Delta}_S^\tau}$$

Notons qu'on introduit ainsi le contrefactuel $\bar{X}_A \hat{\gamma}_B^\tau$, qui correspond à la façon dont les caractéristiques moyennes des individus du groupe A seraient valorisées par les "rendements" que connaissent les B au quantile (non-conditionnel) d'ordre τ . La décomposition détaillée obtenue est bien, tout comme la décomposition d'Oaxaca-Blinder, additive, et indépendante à l'ordre.

Cette méthode est très simple à mettre en œuvre pour les quantiles (surtout sachant que les γ peuvent être directement estimés à l'aide du *package* `uqr`, cf. application 4). L'emploi de la RIF peut également être élargi à d'autres statistiques distributionnelles que les quantiles⁹, notamment au rapport interdécile ou au taux de pauvreté relative. Il faut alors calculer la *RIF* correspondante. Toutefois, on procède à une approximation locale et la qualité de cette approximation pourrait notamment être problématique en présence de points de masse. Cette méthode est ainsi complémentaire de la méthode de repondération de DiNardo et al. (1996), comme le soulignent Fortin et al. (2011) – on peut dans un premier temps appliquer la repondération pour obtenir la décomposition agrégée, puis appliquer les régressions quantiles non-conditionnelles pour parvenir à la décomposition détaillée.

Encadré 5 : La régression quantile non-conditionnelle

Pour estimer les γ , on a recours aux régressions quantiles non-conditionnelles, ou régressions sur RIF (pour *Recentered Influence Function*, ou fonction d'influence recentrée). La fonction d'influence, outil classique en statistiques robustes, appréhende la façon dont une observation particulière Y_i influence une statistique donnée. Dans le cas où la statistique considérée est le quantile d'ordre τ de la distribution de Y (qu'on note Q^τ), la fonction d'influence recentrée associée à Y_i la grandeur suivante : $RIF(Y_i; Q^\tau) = Q^\tau + \frac{\tau - \mathbf{1}\{Y_i \leq Q^\tau\}}{f_Y(Q^\tau)}$. Pour un quantile Q^τ donné, cette fonction ne prendra que deux valeurs selon que Y_i se situe en-dessous ou au-dessus de Q^τ . Si l'on considère par exemple une distribution de salaire dont

9. Pour la moyenne, on retombe sur la régression standard de Y sur X).

la médiane ($\tau = 0.5$) vaut 1700, la fonction d'influence recentrée vaut pour chaque Y_i : $1700 + \frac{0.5 - \mathbf{1}\{Y_i \leq 1700\}}{f_Y(1700)}$.

Une régression quantile non-conditionnelle (au quantile d'ordre τ) correspond ensuite simplement à une régression par MCO de la grandeur $RIF(Y_i; Q^\tau)$ sur X .^a L'obtention des valorisations γ^τ se fait donc à travers deux étapes simples : transformation de chaque Y_i en $RIF(Y_i; Q^\tau)$; puis régression linéaire de $RIF(Y_i; Q^\tau)$ sur les X . En pratique, le package `uqr` permet de procéder directement aux régressions quantiles non-conditionnelles.

^a. Ce faisant, on modélise comme dans une régression linéaire classique $\mathbb{E}([RIF(Y, Q^\tau)|X]) = X \cdot \gamma^\tau + \epsilon$. Or la RIF permet d'écrire $\mathbb{E}[RIF(Y, Q^\tau)] = Q^\tau$, et de là $Q^\tau = \mathbb{E}[RIF(Y, Q^\tau)] = \mathbb{E}_X[\mathbb{E}([RIF(Y, Q^\tau)|X])] = \mathbb{E}[X] \cdot \gamma^\tau$. La contrepartie empirique de cette expression, $\hat{Q}^\tau = \bar{X} \hat{\gamma}^\tau$, permet la décomposition présentée plus haut.

Application 4 : décomposition par régression quantile non-conditionnelle dans R

Le package `uqr` permet d'implémenter des régressions quantiles non-conditionnelles sous R, à travers la fonction `urq`. On spécifie le(s) quantile(s) au(x)quel(s) on souhaite effectuer ces régressions grâce à l'option `tau=`. On effectue cette décomposition séparément pour la population de référence (`ref<-base$magh==0`) d'une part, et pour les descendants d'immigrés (`!ref`) d'autre part.

```
library(uqr)
rif.ref<-urq(formula=logsal~exp+exp2+ddipl, data=base[ref,], tau=grid)

#On obtient par exemple pour les 3 premiers déciles :
rif.ref$coefficients[,1:3]

           tau= 0.1  tau= 0.2  tau= 0.3
(Intercept) 6.909746  7.033857  7.139426
exp          0.027126  0.027002  0.027141
exp2        -0.000417 -0.000388 -0.000363
ddipl1       0.065448  0.089325  0.114753
ddipl4      -0.071736 -0.108879 -0.135830
ddipl5      -0.122215 -0.185484 -0.247785
ddipl6      -0.149501 -0.200277 -0.237959
ddipl7      -0.260060 -0.326616 -0.388492

rif.magh<-urq(formula=logsal~exp+exp2+ddipl, data=base[!ref,], tau=grid)
```

Une fois les valorisations γ obtenues pour chaque groupe, on peut procéder à une décomposition "classique" de type Oaxaca-Blinder, à partir des vecteurs de moyennes calculés pour chaque variable, dans la population de référence et parmi les descendants d'immigrés.

```

#Calcul des X moyens dans chaque groupe
X<-model.matrix(logsal ~ exp+exp2+ddipl1,base)
moy.ref<-apply(X[ref,c(2:8)],2,mean)
moy.magh<-apply(X[!ref,c(2:8)],2,mean)

#Calcul des écarts expliqués par chaque variable
expl.detail<-apply(rif.ref$coefficients[2:8,], 2, "*",moy.ref-moy.magh)
expl.detail[,1:3]

      tau= 0.1 tau= 0.2 tau= 0.3
exp      0.13003  0.12944  0.13011
exp2     -0.09453 -0.08794 -0.08249
ddipl1   0.00121  0.00166  0.00213
ddipl4   0.00260  0.00395  0.00493
ddipl5  -0.00626 -0.00950 -0.01269
ddipl6   0.00172  0.00230  0.00274
ddipl7   0.01271  0.01596  0.01898

```

On obtient alors l'ensemble des contributions détaillées à l'effet de composition, pour chaque variable, en chaque décile. On peut retrouver les parts expliquées et inexpliquées totales en chaque et les comparer aux résultats de la décomposition agrégée obtenus par repondération.

```

rif.expl<-apply(expl.detail,2,sum)

#Ecart expliqués totaux (effet de composition)
round(rif.expl,3)

tau= 0.1 tau= 0.2 tau= 0.3 tau= 0.4 tau= 0.5 tau= 0.6 tau= 0.7 tau= 0.8
      0.047   0.056   0.064   0.070   0.078   0.083   0.090   0.102
tau= 0.9
      0.103

#De même pour les écarts inexpliqués
inexpl.detail<-rbind(rif.ref$coefficients[1,]-rif.magh$coefficients[1,],
                    apply(rif.ref$coefficients[2:8,]-rif.magh$coefficients[2:8,],
                          2, "*",moy.magh))
rif.inexpl<-apply(inexpl.detail,2,sum)

#Ecart inexpliqués totaux
round(rif.inexpl,3)

tau= 0.1 tau= 0.2 tau= 0.3 tau= 0.4 tau= 0.5 tau= 0.6 tau= 0.7 tau= 0.8
      0.021   0.012   0.005   0.006   0.019   0.022   0.043   0.058
tau= 0.9
      0.088

```

L'approximation effectuée par les régressions RIF n'apparaît pas très importante ici, si l'on considère l'écart total mesuré. Par exemple au 9^{ème} décile, la somme de l'écart expliqué et inexpliqué estimés par RIF donne 0.191, tandis que l'écart brut est de 0.188.

Encadré 6 : Les méthodes d'estimation de la distribution conditionnelle

Les régressions sur fonction de répartition

Repartons de la distribution contrefactuelle $F_{Y_B|X_A}$, qui correspond à la distribution conditionnelle des salaires du groupe B appliquée aux caractéristiques des individus du groupe A . Une façon naturelle de l'obtenir est d'estimer directement la distribution conditionnelle $F(Y_B|X)$, distribution des Y comme fonction des caractéristiques X parmi les individus du groupe B , et de l'appliquer aux caractéristiques du groupe A .

La fonction de distribution peut être résumée par un ensemble de probabilités de se situer en-dessous d'un certain seuil : pour reprendre l'exemple des salaires, on souhaiterait modéliser le fait d'avoir un salaire inférieur à 1500 € par mois, d'avoir un salaire inférieur à 2000 € par mois, etc., en raisonnant à chaque fois à caractéristiques données. Le problème revient ainsi à réaliser un ensemble d'estimations sur des indicatrices (se situer au-dessus ou en-dessous d'un seuil donné), ce qui peut être fait de façon très classique par un modèle logit, un modèle probit, ou même un modèle de probabilité linéaire. Plus le nombre de seuils considérés sera grand, plus l'estimation de la fonction de répartition sera fine. C'est sur cette idée que repose la méthode de régression sur distribution proposée par Chernozhukov et al. (2013).

En pratique, la méthode de Chernozhukov et al. (2013) consiste donc en trois étapes, pour chaque seuil $y \in [\min(Y), \max(Y)]$:

- (1) On estime au sein du groupe A un logit, un probit ou une régression linéaire sur $\mathbf{1}\{Y_i \leq y\}$, avec $P(Y < y|X) = F(X\beta(y))$, selon que $F(\cdot)$ est la fonction de répartition de la loi logistique, de la loi normale, ou la fonction identité.
- (2) On utilise les coefficients estimés à l'étape (1) pour calculer la probabilité prédite $F(X_i\hat{\beta}_A(y))$ pour chaque individu i du groupe B .
- (3) On calcule la moyenne de ces probabilités prédites sur l'ensemble du groupe B pour finalement obtenir $\hat{F}_{Y_A^C}(y) = \frac{1}{N_B} \sum_{i \in B} F(X_i\hat{\beta}_A(y))$, c'est-à-dire la probabilité de se situer au-dessous du seuil y qu'auraient les individus du groupe B si leurs caractéristiques étaient valorisées de la même façon que pour les A .

On obtient ainsi un ensemble de probabilités contrefactuelles qui permettent de reconstituer la distribution souhaitée. Il est toutefois possible que la fonction de distribution estimée ne soit pas monotone : pour y_1 et y_2 proches avec $y_1 < y_2$, rien ne garantit que $\hat{F}_{Y_A^C}(y_1) < \hat{F}_{Y_A^C}(y_2)$. Il est alors nécessaire d'utiliser une procédure de lissage pour s'assurer qu'elle pourra bien être inversée en quantiles.

Cette méthode est intensive en calculs : elle demande d’effectuer un grand nombre de régressions lorsqu’on souhaite inverser la distribution et revenir à l’écart entre quantiles. On peut toutefois souhaiter simplement décomposer l’écart de probabilité de se trouver au-dessus ou en-dessous d’un certain niveau de salaire, c’est même parfois plus parlant (exemple : seuil de pauvreté en termes absolus ; “hauts revenus” comparable dans les deux groupes). On a alors besoin d’effectuer la procédure uniquement pour le y d’intérêt (ce qui correspond à la méthode de Fairlie, voir section 3.2). En cela, il est plus simple d’estimer des distances “verticales” que des distances “horizontales” (Figure ??). Notons que comme pour Fairlie, si de plus on modélise cette probabilité par un modèle de probabilité linéaire, on retombe sur une décomposition de type Oaxaca-Blinder.

Estimation de distributions contrefactuelles par régressions quantiles

Il existe d’autres manières que celle proposée par Chernozhukov et al. (2013) de construire une distribution contrefactuelle du groupe A en mimant la distribution conditionnelle du groupe B.

Machado and Mata (2005) et Melly (2005) utilisent ainsi un procédé de transformation de chaque Y_{Ai} en un $Y_{Bi: X=X|D_A}^C$ via la simulation d’un ensemble de quantiles.

- (1) On tire un ensemble $\tau_1, \tau_2, \dots, \tau_S$ entre 0 et 1.
- (2) En chacun de ces quantiles τ_s on estime une régression quantile parmi le groupe B (cela permet d’estimer une fonction de rendement des caractéristiques au sein de ce groupe, en S points de la distribution).
- (3) Les rendements estimés chez les B permet de prédire un Y_{As}^C à partir des X de chaque individu du groupe A.

En parcourant chaque quantile de la distribution conditionnelle des B et en intégrant sur les X des A, on retrouve la distribution contrefactuelle d’intérêt. L’avantage de ce procédé par rapport à celui de Juhn et al. (1993) est qu’il permet de tenir compte de l’hétéroscédasticité des résidus.

Cette méthode a aussi cependant des limites : elle impose de faire la simulation en un grand nombre de points, même dans le cas où l’on ne s’intéresse qu’à la décomposition en un seul quantile de la distribution. La procédure est ainsi très intensive en calculs (même si la version de Melly (2005) consistant à tirer à chaque itération un ensemble de X dans l’échantillon des A l’est un peu moins). De plus, la spécification linéaire peut être restrictive dans certaines applications, notamment dans le cas où la distribution des Y présente des points de masse (cela peut être le cas pour une distribution de salaire en présence d’un salaire minimum).

Pour l’ensemble de ces méthodes d’estimation de la distribution conditionnelle, la décomposition détaillée est possible mais là encore elle ne peut être à la fois additive et indépendante à l’ordre des variables.

Références

- John Abowd, Francis Kramarz, and David Margolis. High wage workers and high wage firms. *Econometrica*, 67(2) :251–333, 1999.
- Joseph G. Altonji, Prashant Bharadwaj, and Fabian Lange. Changes in the Characteristics of American Youth : Implications for Adult Outcomes. *Journal of Labor Economics*, 30(4) :783 – 828, 2012.
- David Audenaert, José Bardaji, Raphaël Lardeux, Michaël Orand, and Michaël Sicsic. La résistance des salaires depuis la grande récession s’explique-t-elle par des rigidités à la baisse ? *Insee Références, L’économie française - Comptes et dossiers*, 2014.
- Christophe Bertran. Le revenu d’activité des non-salariés : plus élevé en moyenne dans les départements du nord que dans ceux du sud. *Insee Première*, (1672), 2017.
- Alan Blinder. Wage discrimination : reduced form and structural estimates. *Journal of Human resources*, (1672), 1973.
- David Card, Ana Cardoso, and Patrick Kline. Bargaining, sorting, and the gender wage gap : Quantifying the impact of firms on the relative pay of women. *The Quarterly Journal of Economics*, 131(2) :633–686, 2016.
- Victor Chernozhukov, Iván Fernández-Val, and Blaise Melly. Inference on counterfactual distributions. *Econometrica*, 81(6) :2205–2268, 2013.
- Jeremiah Cotton. On the decomposition of wage differentials. *The review of economics and statistics*, pages 236–243, 1988.
- John DiNardo, Nicole Fortin, and Thomas Lemieux. Labour market institutions and the distribution of wages, 1973-1992 : a semi parametric approach. *Econometrica*, 64(5) :1001, 1996.
- Robert W. Fairlie. An extension of the blinder-oaxaca decomposition technique to logit and probit models. *Journal of economic and social measurement*, 30(4) :305–316, 2005.
- Sergio Firpo, Nicole M. Fortin, and Thomas Lemieux. Decomposing Wage Distributions using Recentered Influence Functions Regressions. mimeo, University of British Columbia, 2007.
- Sergio Firpo, Nicole M. Fortin, and Thomas Lemieux. Unconditional Quantile Regressions. *Econometrica*, 77(3) :953–973, 05 2009.
- Nicole Fortin, Thomas Lemieux, and Sergio Firpo. Decomposition methods in economics. *Handbook of labor economics*, 4 :1–102, 2011.
- Javier Gardeazabal and Arantza Ugidos. More on identification in detailed wage decompositions. *Review of Economics and Statistics*, 86(4) :1034–1036, 2004.
- Keisuke Hirano, Guido W. Imbens, and Geert Ridder. Efficient Estimation of Average Treatment Effects Using the Estimated Propensity Score. *Econometrica*, 71(4) :1161–1189, 07 2003.

- Marek Hlavac. *oaxaca* : Blinder-oaxaca decomposition in r. 2014.
- Ben Jann. The blinder-oaxaca decomposition for linear regression models. *The Stata Journal*, 8(4) :453–479, 2008.
- Frank L. Jones and Jonathan Kelley. Decomposing differences between groups a cautionary note on measuring discrimination. *Sociological Methods & Research*, 12(3) :323–343, 1984.
- Chinhui Juhn, Kevin M. Murphy, and Brooks Pierce. Wage inequality and the rise in returns to skill. *Journal of political Economy*, pages 410–442, 1993.
- Claire Kubrak. Principe et mise en oeuvre des approches comptable et économétrique. *Document de travail Insee-Direction de la Diffusion et de l'Action régionale*, H 2018/01, 2018.
- Rasmus Lentz and Dale T. Mortensen. Labor market models of worker and firm heterogeneity. *Annual Review of Economics*, 2(1) :577–602, 2010.
- José Machado and José Mata. Counterfactual decomposition of changes in wage distributions using quantile regression. *Journal of applied Econometrics*, 20(4) :445–465, 2005.
- Blaise Melly. Decomposition of differences in distribution using quantile regression. *Labour economics*, 12(4) :577–590, 2005.
- David Neumark. Employers' discriminatory behavior and the estimation of wage discrimination. *Journal of Human Resources*, 23(3) :279–295, 1988.
- Ronald Oaxaca. Male-female wage differentials in urban labor markets. *International Economic Review*, pages 693–709, 1973.
- Ronald Oaxaca and Michael Ransom. Identification in detailed wage decompositions. *Review of Economics and Statistics*, 81(1) :154–157, 1999.
- Cordelia Reimers. Labor market discrimination against hispanic and black men. *The review of economics and statistics*, pages 570–579, 1983.
- Myeong-Su Yun. Decomposing differences in the first moment. *Economics letters*, 82(2) :275–280, 2004.
- Myeong-Su Yun. A simple solution to the identification problem in detailed wage decompositions. *Economic inquiry*, 43(4) :766–772, 2005.
- Myeong-Su Yun. Identification problem and detailed oaxaca decomposition : A general solution and inference. *Journal of economic and social measurement*, 33(1) :27–38, 2008.