

Estimating a counterfactual wage distribution using survey data

Mihaela-Cătălina Anastasiade¹, Alina Matei² and Yves Tillé²

Swiss Federal Statistical Office¹, University of Neuchâtel², Switzerland

Colloque francophone sur les sondages 2018
Lyon

- We present a parametric approach to estimate 'wage discrimination' at different quantiles.
- The goal is to reduce the variance of the estimates compared to the existing methods.
- We illustrate this approach using the generalized beta of the second kind distribution (hereafter, GB2).

- Consider a finite population with the labels $U = \{1, 2, \dots, N\}$.
- From this population, we randomly select a sample S of size n , without replacement.
- The sample is selected through a sampling design $p(s) = \Pr(S = s), \forall s \subseteq U$.
- To each unit $k \in S$, a survey weight w_k is associated.
- These weights can be equal to the inverse of the inclusion probabilities or can be more complicated weights, like calibration weights.
- Let y be the variable of interest (the wage).

Superpopulation framework

- Assume that $U = U_M \cup U_F$, $U_M \cap U_F = \emptyset$ is drawn from a superpopulation.
- The superpopulation is also divided in two subsuperpopulations from where the subsets U_g , $g = \{M, F\}$ are drawn, respectively.
- The wage is a random variable Y_g and \mathbf{X}_g is a set of covariates.
- In each subset U_g ,

$$Y_{k,g} \mid \mathbf{X}_g = \mathbf{x}_{k,g} \sim D(\gamma_{k,g}, \boldsymbol{\delta}_g), k \in U_g.$$

- We assume that $D(\gamma_{k,g}, \boldsymbol{\delta}_g)$, $k \in U_g$ is a continuous distribution and that $\gamma_{k,g} = h(\mathbf{x}_{k,g}^\top \boldsymbol{\beta}_g)$, where h is a known continuous function.
- The sample $S = S_M \cup S_F$, $S_M \cap S_F = \emptyset$, $S_g \subseteq U_g$, $g = \{M, F\}$.

The counterfactual wage distribution

The CDF of the counterfactual wage distribution is defined as

$$F^C(y) = \int_{\mathcal{X}_M} F^{Y_F|\mathbf{X}_F}(y | \mathbf{x}) dF^{\mathbf{X}_M}(\mathbf{x}),$$

where \mathcal{X}_M is the support of \mathbf{X}_M and \mathcal{X}_F the support of \mathbf{X}_F .

It is assumed that $\mathcal{X}_M \subseteq \mathcal{X}_F$.

It is interpreted as the distribution function of wages that would be obtained for women if their characteristics were same as those of men.

$$F^F(y) = \int_{\mathcal{X}_F} F^{Y_F|\mathbf{X}_F}(y | \mathbf{x}) dF^{\mathbf{X}_F}(\mathbf{x}),$$

$$F^M(y) = \int_{\mathcal{X}_M} F^{Y_M|\mathbf{X}_M}(y | \mathbf{x}) dF^{\mathbf{X}_M}(\mathbf{x}).$$

The counterfactual wage distribution

The CDF of the counterfactual wage distribution is defined as

$$\begin{aligned} F^C(y) &= \int_{\mathcal{X}_M} F^{Y_F | \mathbf{X}_F}(y | \mathbf{x}) dF^{\mathbf{X}_M}(\mathbf{x}) \\ &= \int_{\mathcal{X}_F} F^{Y_F | \mathbf{X}_F}(y | \mathbf{x}) \frac{dF^{\mathbf{X}_M}(\mathbf{x})}{dF^{\mathbf{X}_F}(\mathbf{x})} dF^{\mathbf{X}_F}(\mathbf{x}) \\ &= \int_{\mathcal{X}_F} F^{Y_F | \mathbf{X}_F}(y | \mathbf{x}) \psi(\mathbf{x}) dF^{\mathbf{X}_F}(\mathbf{x}), \end{aligned}$$

where \mathcal{X}_M is the support of \mathbf{X}_M and \mathcal{X}_F the support of \mathbf{X}_F . It is assumed that $\mathcal{X}_M = \mathcal{X}_F$.

$$F^F(y) = \int_{\mathcal{X}_F} F^{Y_F | \mathbf{X}_F}(y | \mathbf{x}) dF^{\mathbf{X}_F}(\mathbf{x}),$$

The weighted DiNardo, Fortin and Lemieux method

- DiNardo et al. (1996) write the reweighting factor $\psi(\mathbf{x}_k) = \frac{dF^{\mathbf{X}_M}(\mathbf{x}_k)}{dF^{\mathbf{X}_F}(\mathbf{x}_k)}$ as

$$\psi(\mathbf{x}_k) = \psi_k = \frac{P(\text{Gender}_k = \text{'man'} \mid \mathbf{x}_k) / P(\text{Gender}_k = \text{'man'})}{P(\text{Gender}_k = \text{'woman'} \mid \mathbf{x}_k) / P(\text{Gender}_k = \text{'woman'})}$$

- The idea is to reweigh the characteristics of women so that they match the characteristics of men, such that

$$\widehat{\mathbf{X}}_C = \widehat{\mathbf{X}}_M,$$

where $\widehat{\mathbf{X}}_C = \sum_{k \in S_F} \widehat{\psi}_k w_k \mathbf{x}_k / \sum_{k \in S_F} \widehat{\psi}_k w_k$ and

$$\widehat{\mathbf{X}}_M = \sum_{k \in S_M} w_k \mathbf{x}_k / \sum_{k \in S_M} w_k.$$

- The factor $\psi(\mathbf{x}_k)$ can be estimated by using a probit or a logistic regression model (DiNardo et al., 1996) or by calibration (Anastasiade and Tillé, 2017).

If $\Delta_{(\alpha)}$ is the wage difference between men and women at a given quantile α , we can write

$$\Delta_{(\alpha)} = Q_{(\alpha)}^M - Q_{(\alpha)}^F = (Q_{(\alpha)}^M - Q_{(\alpha)}^C) + (Q_{(\alpha)}^C - Q_{(\alpha)}^F),$$

where $Q_{(\alpha)}^M$, $Q_{(\alpha)}^C$ and $Q_{(\alpha)}^F$ are the quantile of order α of men, counterfactual and women wage distributions, respectively.

$$\hat{\Delta}_{(\alpha)} = \hat{Q}_{(\alpha)}^M - \hat{Q}_{(\alpha)}^F = (\hat{Q}_{(\alpha)}^M - \hat{Q}_{(\alpha)}^C) + (\hat{Q}_{(\alpha)}^C - \hat{Q}_{(\alpha)}^F).$$

Estimation of the CDF in finite populations

- The empirical CDF at the U level is defined as

$$F_{emp}(y) = \frac{\sum_{k \in U} I(y_k \leq y)}{N}.$$

- The classical design-based estimator is

$$\hat{F}_{emp}(y) = \frac{\sum_{k \in S} w_k I(y_k \leq y)}{\sum_{k \in S} w_k}.$$

- The quantile of order α of y is estimated by

$$\hat{Q}_{\alpha, emp}(y) = \inf \{ \hat{F}_{emp}(y) \geq \alpha \}.$$

Parametric approach

- We assume that U is selected from a superpopulation and Y is a random variable with the CDF $F(\cdot)$.
- We include auxiliary information \mathbf{X} in the estimation of $F(\cdot)$ by assuming that $Y_k | \mathbf{X}_k \sim D(\gamma_k = h(\mathbf{x}_k^\top \boldsymbol{\beta}), \boldsymbol{\delta}_k), k \in U$.
- We write

$$F_U(y) = \sum_{k \in U} \lambda_k F_{D(\gamma_k, \boldsymbol{\delta}_k)}(y | \mathbf{x}_{k,g}) = \frac{1}{N} \sum_{k \in U} F_{D(\gamma_k, \boldsymbol{\delta}_k)}(y | \mathbf{x}_{k,g}),$$

where $\lambda_k = 1/N$, $F_{D(\gamma_k, \boldsymbol{\delta}_k)}(\cdot | \mathbf{x}_k)$ is the CDF of the distribution $D(\gamma_k = h(\mathbf{x}_k^\top \boldsymbol{\beta}), \boldsymbol{\delta}_k), k \in U$, and $h(\cdot)$ is a continuous function.

Method 1 for quantile estimation

We propose to estimate the quantile of order α of Y as

$$\hat{Q}_{(\alpha)} = \inf\{y \mid \hat{F}_U(y) \geq \alpha\},$$

where $\hat{F}_U(y)$ is the estimator of $F_U(y)$ in the point y given by

$$\hat{F}_U(y) = \sum_{k \in S} w_k \hat{F}_{D(\hat{\gamma}_k, \hat{\delta})}(y_k \mid \mathbf{x}_k) / \sum_{k \in S} w_k.$$

Method 2 for quantile estimation

In case the inverse function of $\hat{F}_U(y)$ cannot be computed, we propose to use a Monte Carlo method based on parametric bootstrap.

Method 2 for quantile estimation

$$Y_{i,k} | \mathbf{x}_k \sim D(h(\mathbf{x}_k^\top \hat{\boldsymbol{\beta}}), \hat{\boldsymbol{\delta}})$$

$$\begin{pmatrix} w_1 & w_2 & w_3 & \dots & w_n \\ y_{11} & y_{12} & y_{13} & \dots & y_{1n} \\ y_{21} & y_{22} & y_{23} & \dots & y_{2n} \\ y_{31} & y_{32} & y_{33} & \dots & y_{3n} \\ \dots & \dots & \dots & \dots & \dots \\ y_{m1} & y_{m2} & y_{m3} & \dots & y_{mn} \end{pmatrix} \longrightarrow \begin{pmatrix} \hat{Q}_1^{0.25} & \hat{Q}_1^{0.5} & \hat{Q}_1^{0.75} & \hat{Q}_1^{0.95} \\ \hat{Q}_2^{0.25} & \hat{Q}_2^{0.5} & \hat{Q}_2^{0.75} & \hat{Q}_2^{0.95} \\ \hat{Q}_3^{0.25} & \hat{Q}_3^{0.5} & \hat{Q}_3^{0.75} & \hat{Q}_3^{0.95} \\ \dots & \dots & \dots & \dots \\ \hat{Q}_m^{0.25} & \hat{Q}_m^{0.5} & \hat{Q}_m^{0.75} & \hat{Q}_m^{0.95} \\ \hat{Q}_{0.25} & \hat{Q}_{0.5} & \hat{Q}_{0.75} & \hat{Q}_{0.95} \end{pmatrix}$$

Remark: Methods 1 and 2 are applied to estimate respectively the α -quantiles in each group $g \in \{M, F\}$.

Method 1 for the counterfactual wage distribution

We redefine the counterfactual CDF at the U_F level as

$$F_{U_F}^C(y) = \frac{1}{N_C} \sum_{k \in U_F} \psi_k F^{(Y_F | \mathbf{x}_F)}(y_k | \mathbf{x}_{k,F}),$$

where $N_C = \sum_{k \in U_F} \psi_k$.

- First, we estimate it by

$$\hat{F}_{U_F}^C(y) = \frac{\sum_{k \in S_F} \hat{\psi}_k w_k \hat{F}^{(Y_F | \mathbf{x}_F)}(y_k | \mathbf{x}_{k,F})}{\sum_{k \in S_F} \hat{\psi}_k w_k},$$

where $\hat{F}^{(Y_F | \mathbf{x}_F)}(y_k | \mathbf{x}_{k,F}) = F_{D(h(\mathbf{x}_{k,F}^\top \hat{\beta}_F), \hat{\delta}_F)}(y_k | \mathbf{x}_{k,F})$, and $\hat{\psi}_k$ is estimated by calibration (Anastasiade and Tillé, 2017).

- Next, $\hat{Q}_{(\alpha)}^C = \inf\{y | \hat{F}_{U_F}^C(y) \geq \alpha\}$.

Method 2 for the counterfactual wage distribution

$$Y_{i,k} \mid \mathbf{x}_{k,F} \sim D(h(\mathbf{x}_{k,F}^\top \widehat{\boldsymbol{\beta}}_F), \widehat{\boldsymbol{\delta}}_F)$$

$$\left(\begin{array}{cccc} \widehat{\psi}_1 w_1 & \widehat{\psi}_2 w_2 & \dots & \widehat{\psi}_{n_F} w_{n_F} \\ y_{11} & y_{12} & \dots & y_{1n_F} \\ y_{21} & y_{22} & \dots & y_{2n_F} \\ y_{31} & y_{32} & \dots & y_{3n_F} \\ \dots & \dots & \dots & \dots \\ y_{m1} & y_{m2} & \dots & y_{mn_F} \end{array} \right) \rightarrow \left(\begin{array}{cccc} \widehat{Q}_1^{0.25} & \widehat{Q}_1^{0.5} & \widehat{Q}_1^{0.75} & \widehat{Q}_1^{0.95} \\ \widehat{Q}_2^{0.25} & \widehat{Q}_2^{0.5} & \widehat{Q}_2^{0.75} & \widehat{Q}_2^{0.95} \\ \widehat{Q}_3^{0.25} & \widehat{Q}_3^{0.5} & \widehat{Q}_3^{0.75} & \widehat{Q}_3^{0.95} \\ \dots & \dots & \dots & \dots \\ \widehat{Q}_m^{0.25} & \widehat{Q}_m^{0.5} & \widehat{Q}_m^{0.75} & \widehat{Q}_m^{0.95} \\ \widehat{Q}_{0.25}^C & \widehat{Q}_{0.5}^C & \widehat{Q}_{0.75}^C & \widehat{Q}_{0.95}^C \end{array} \right)$$

The two methods

- The proposed methods aim to reduce the variance of the estimated quantiles compared to the estimation given by the empirical CDF.
- The methods are correct if the conditional distribution is correct.
- Departures from this assumption can be managed by using a GB2 distribution.

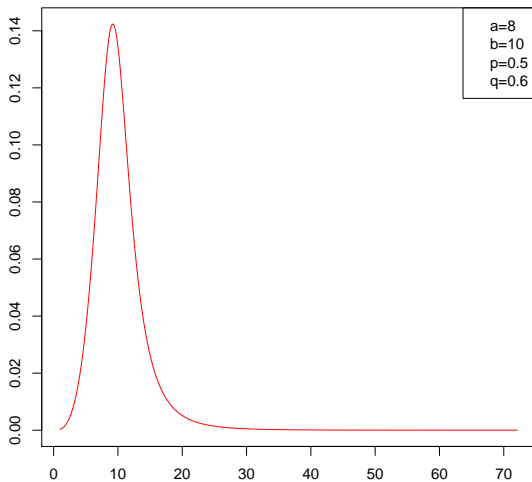
The GB2 distribution

The GB2 distribution is characterized by four parameters, namely a , b , p and q . The probability density function of a $GB2(a, b, p, q)$ distribution is given by

$$f(y; a, b, p, q) = \frac{a \left(\frac{y}{b}\right)^{ap-1}}{bB(p, q) \left[1 + \left(\frac{y}{b}\right)^a\right]^{p+q}},$$

where a , p , q are the shape parameters and b is a scale parameter.

Example of GB2 distribution

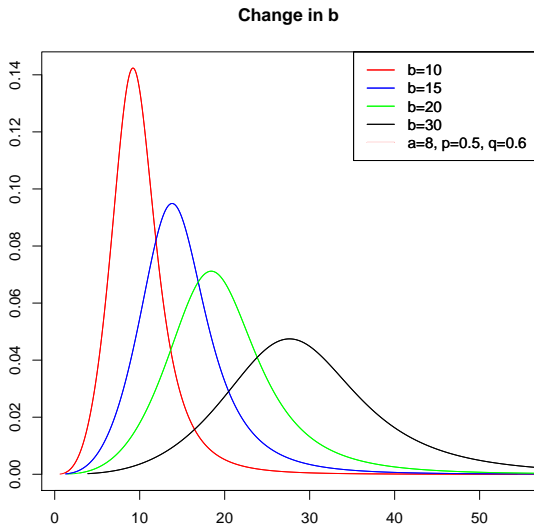


The GB2 distribution

- We borrow from McDonald and Butler (1990) the idea of changing the scale parameter, by expressing it as a function of the observed characteristics.
- In each group, $g \in \{M, F\}$, we assume that the conditional wage of $k \in U_g$, $Y_k \mid \mathbf{X}_{k,g} = \mathbf{x}_{k,g} \sim GB2(a_g, \exp(\mathbf{x}_k \boldsymbol{\beta}_g), p_g, q_g)$.
- Thus, for each $k \in U_g$, the GB2 density becomes

$$f[y_k; a_g, \exp(\mathbf{x}_k^\top \boldsymbol{\beta}_g), p_g, q_g] = \frac{a \left[\frac{y_k}{\exp(\mathbf{x}_k^\top \boldsymbol{\beta}_g)} \right]^{a p_g - 1}}{\exp(\mathbf{x}_k^\top \boldsymbol{\beta}_g) B(p_g, q_g) \left\{ 1 + \left[\frac{y_k}{\exp(\mathbf{x}_k^\top \boldsymbol{\beta}_g)} \right]^a \right\}^{p_g + q_g}}.$$

The GB2 distribution



The GB2 regression

We assume that

$$\log(Y_{k,g}) = \mathbf{X}_{k,g}^\top \boldsymbol{\beta}_g + \log(\varepsilon_{k,g}),$$

where $Y_{k,g}$ is the wage of individual $k \in U_g$, $\varepsilon_{k,g} \sim GB2(a_g, 1, p_g, q_g)$.

The GB2 distribution

- We estimate the parameters of the GB2 distribution using pseudo-maximum likelihood.
- We developed an algorithm to estimate the parameters of the GB2 distribution when x_k is introduced in the scale parameter.
- We estimate the standard errors of the estimated parameters using the sandwich estimator and a parametric bootstrap approach.

Example - Swiss Survey on Earnings Structure, 2012

- sample of 5643 employees (1880 women, 3763 men) working in the economic activity 'Manufacture of computer, electronic and optical products';
- models with the covariates: age, education level (9 categories), professional position (5 categories).

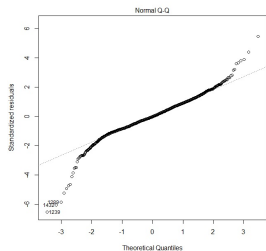


Figure: QQplot for a log-normal model

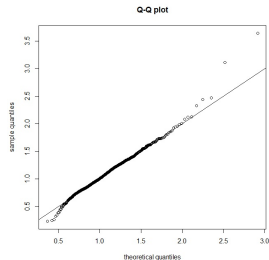
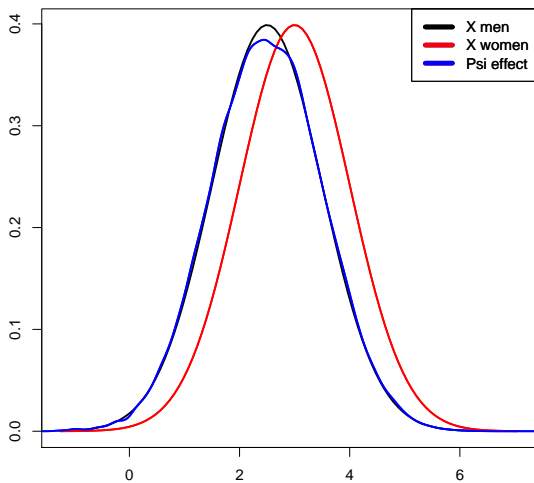


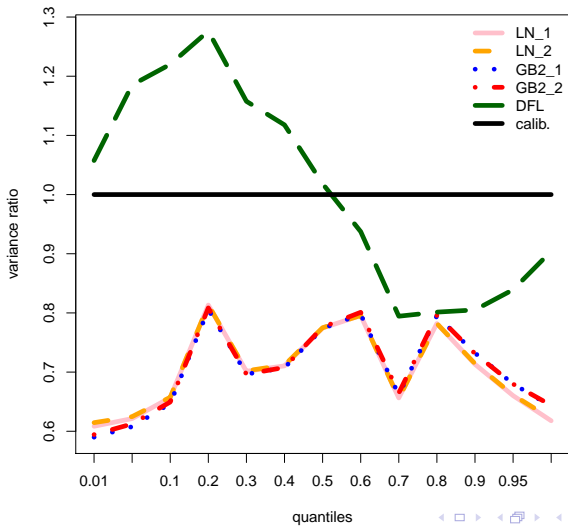
Figure: QQplot for a GB2 model

- $N_F = N_M = 50,000, n_F = n_M = 10,000,$
- $X_{k,F} \sim N(3, 1), X_{k,M} \sim N(2.5, 1),$ independent,
- $Y_{k,F} \sim LN(1.15 + 2.5X_{k,F}, 1), k \in U_F$
- 1000 independent srswor samples of size n_F from $U_F,$
- At the population level: $\psi_k = f_{N(2.5,1)}(X_{k,F})/f_{N(3,1)}(X_{k,F}), k \in U_F.$

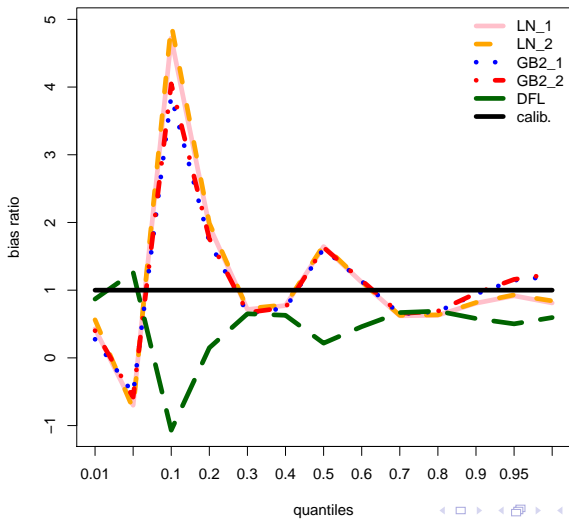
ψ effect



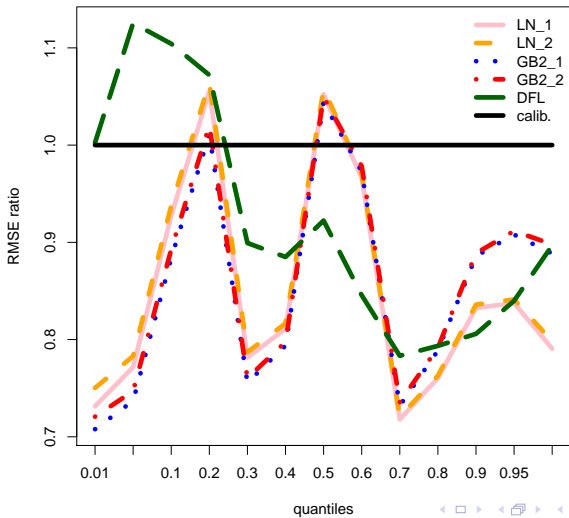
Monte-Carlo variance



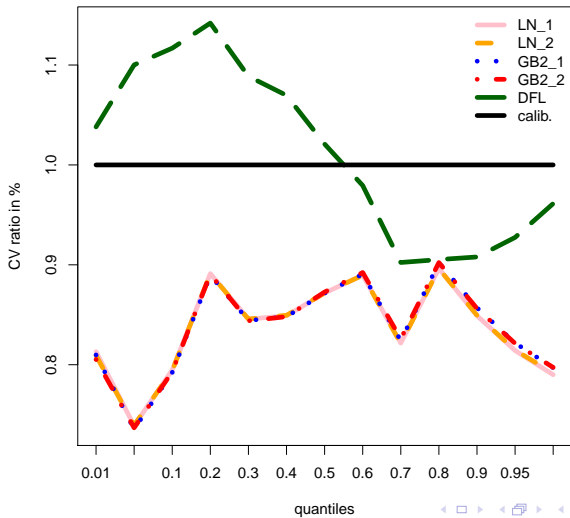
Monte-Carlo bias



Monte-Carlo RMSE



Monte-Carlo CV



- Wages usually have heavy-tailed distributions, which makes it difficult to fit a distribution for them.
- We propose two parametric methods to estimate the quantiles (and differences between the quantiles).
- The introduction of the covariates aims to reduce the variance of the estimates.

References

- Anastasiade, M.-C. and Tillé, Y. (2017). Decomposition of gender wage inequalities through calibration: Application to the swiss structure of earnings survey. *Survey Methodology*, 43(2):211–234.
- DiNardo, J., Fortin, N. M., and Lemieux, T. (1996). Labor Market Institutions and the Distribution of Wages, 1973-1992: A Semiparametric Approach. *Econometrica*, 64(5):1001–44.
- McDonald, J. B. and Butler, R. J. (1990). Regression models for positive random variables. *Journal of Econometrics*, 43(1-2):227–251.