

Generic inference on quantile and quantile effect functions for discrete outcomes

Victor Chernozhukov
Ivan Fernandez-Val
Blaise Melly
Kaspar Wüthrich

The Institute for Fiscal Studies
Department of Economics, UCL

cemmap working paper CWP23/17

Generic Inference on Quantile and Quantile Effect Functions for Discrete Outcomes

Victor Chernozhukov, MIT*
Ivan Fernandez-Val, Boston University
Blaise Melly, University of Bern
Kaspar Wüthrich, UC San Diego

April 6, 2017

Abstract

This paper provides a method to construct simultaneous confidence bands for quantile and quantile effect functions for possibly discrete or mixed discrete-continuous random variables. The construction is generic and does not depend on the nature of the underlying problem. It works in conjunction with parametric, semiparametric, and nonparametric modeling strategies and does not depend on the sampling schemes. It is based upon projection of simultaneous confidence bands for distribution functions. We apply our method to analyze the distributional impact of insurance coverage on health care utilization and to provide a distributional decomposition of the racial test score gap. Our analysis generates new interesting findings, and complements previous analyses that focused on mean effects only. In both applications, the outcomes of interest are discrete rendering standard inference methods invalid for obtaining uniform confidence bands for quantile and quantile effects functions.

Keywords: treatment effects, distribution, discrete, count data, confidence bands, uniform inference.

*We would like to thank for useful comments and feedback participants at various seminars and the students of the courses 14.382 at MIT and EC709 at Boston University, where the ideas presented here have been taught for several years. We would like to acknowledge the financial support from the NSF and from the Swiss National Science Foundation for the grant 165621.

1 Introduction

The quantile function (QF), introduced by Galton (1874), has become a standard tool for descriptive and inferential analysis due to its straightforward and intuitive interpretation. Doksum (1974) suggested to report the quantile effect (QE) function – the difference between two QFs – to compare the distribution of an outcome between two different populations. For example, in Section 4.4, we analyze the racial test score gap by taking the difference of the QFs of the IQ test scores between white and black children. Looking at this QE function allows us to describe the gap not only at the center, but also at the tails of the test score distribution. In randomized control trials and natural experiments, QEs have a causal interpretation, and are usually referred to as quantile treatment effects (QTEs). In Section 4.3, we estimate the treatment effect of insurance coverage on health care utilization based on a conditionally randomized experiment.

Methods for inference on QFs and QE functions are well established when the outcome has a continuous distribution. Under appropriate regularity conditions, properly rescaled and centered empirical analogs of the QFs and QE functions converge to Gaussian processes. In many interesting applications, however, the outcome is not continuously distributed. This is naturally the case for count data, ordinal data, and discrete duration data, but it also concerns test scores that are functions of a finite number of questions, censored variables, and other mixed discrete-continuous variables. Examples include the number of doctor visits in our first application (see Panel A in Figure 1), IQ test scores for children in our second application (see Panels B and C in Figure 1), and wages that have mass points at round values and at the minimum wage. The QFs and QE functions are still well-defined and preserve their intuitive appeal for these types of outcomes, but the *existing* methods are *not* suitable for making inference on them.

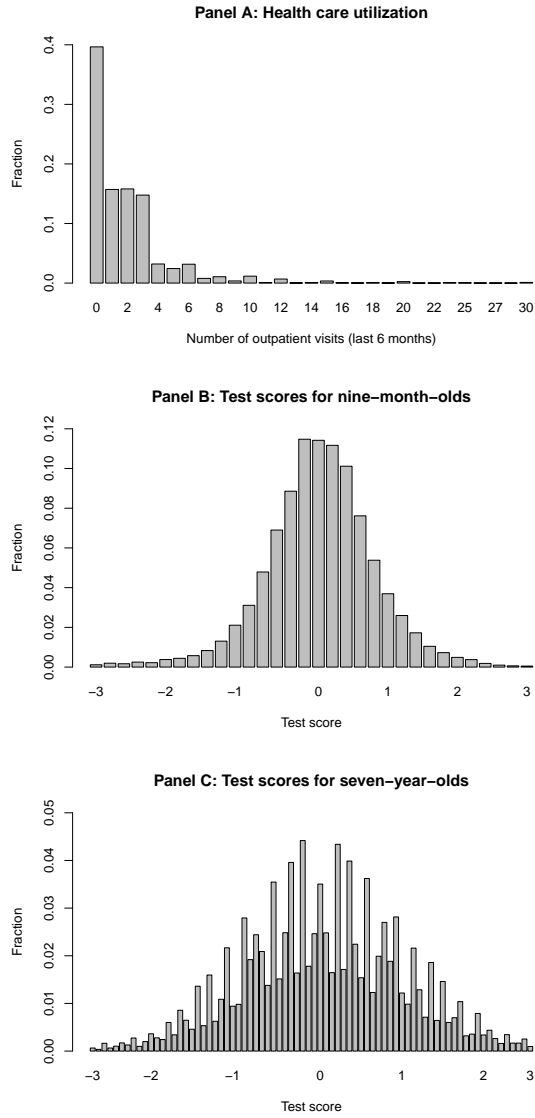


Figure 1: Histograms of the outcomes in our empirical examples. Panel A shows the outcome of our first application reported in Section 4.3; Panel B and C show the outcomes of our second application reported in Section 4.4. Each unique value of the variables has been assigned its own bin.

When the outcome is a discrete or mixed random variable, the empirical quantiles at some probability indexes – those at which the QF jumps – are not even consistent for the corresponding population quantiles, while at other probability indexes – those at which the QF is flat – the empirical quantiles converge to the population quantiles at a rate faster than $1/\sqrt{n}$. Standard inference procedures based on asymptotic Gaussianity cannot perform satisfactorily in such a set-up.

In this paper, we provide a generic construction of *simultaneous* confidence bands for three types of important functions: (1) distribution functions (DFs), (2) QFs, and (3) QE functions. Here simultaneity not only means that the bands are uniform – in that they cover the whole function – but also that all listed objects are covered by the corresponding bands jointly with the prescribed probability. We construct confidence bands for QFs and QE functions from simultaneous confidence bands for DFs. We use *inversion* and *shape* and *support* restrictions to go from confidence bands for two DFs to confidence bands for two QFs and then we take the Minkowski difference of the confidence bands for the QFs, viewed as sets, to construct a confidence band for the QE function.

Our construction is generic and does not depend on the nature of the underlying problem. It applies to the canonical empirical DF, but also works in conjunction with the classical and modern parametric, semiparametric, and nonparametric modeling and estimation strategies, and does not depend on the sampling scheme. For example, it may be employed in conjunction with Poisson, negative binomial, and zero-inflated versions of these models, estimated by the maximum likelihood methods; or it can be employed in conjunction with more flexible methods such as distribution regression, discussed below. Moreover, the QF of interest may be conditional or marginal, counterfactual, or derived from a structural model (see, e.g., Chernozhukov et al. (2013), Imbens and Newey (2009)). The only requirement is the existence of a valid method to obtain simultaneous confidence

bands for DFs. To implement our method, we provide explicit algorithms based on generic bootstrappable estimators of the DFs. By construction, our bands jointly cover all true functions of interest – distribution, quantile, and QE functions – if and only if the bands for the DFs cover the true DFs. Thus, they are not conservative when the bands for the DFs are not conservative, which is asymptotically the case in our – and many other – applications.

To the best of our knowledge this is the first paper that provides simultaneous confidence bands for the whole QFs and QE functions of possibly discrete outcomes. Scheffe and Tukey (1945) were the first to consider inference on quantiles of possibly discrete outcomes, based upon the empirical DF. They show that pointwise confidence intervals for quantiles obtained by projection of the pointwise confidence intervals for the DF are still valid but conservative (asymptotically non-similar) when the outcome is discrete. Frydman and Simon (2008) and Larocque and Randles (2008) suggested methods to estimate the exact coverage rate of these confidence intervals. Our confidence bands for the QFs are uniform and not conservative; in addition, we provide confidence bands for the QE functions, which were not considered before for possibly discrete outcomes.

Another strand of the literature tried to overcome the discreteness in the data by adding a small random noise to the outcome (also called jittering), see for instance Machado and Silva (2005) and the applications in Koenker and Xiao (2002) and Chernozhukov et al. (2013). Ma et al. (2011) considered an alternative definition of quantiles based on linearly interpolated DFs. These strategies restore asymptotic Gaussianity of the empirical QFs and QE functions, at the price of changing the estimand. One might argue that this change is not a serious issue when the number of points in the support of the outcome is large, but we find more transparent to work directly with the observed discrete outcome. Moreover, in a Monte-Carlo study calibrated to a wage distribution, Chernozhukov et al. (2013) provide

evidence supporting the use of approaches that respect the discreteness of outcome, even in cases where the number of support points is not small.

One important application of our generic results concerns the estimation of QFs and QE functions with covariates. When the outcome is continuous, the quantile regression (QR) method, introduced by Koenker and Bassett Jr (1978), is convenient to incorporate covariates. For discrete outcomes, however, the existing inference methods for QR break down. In addition, the linearity assumption for the conditional quantile function underlying QR is highly implausible in that case. For instance, a Poisson regression model does not have linear conditional quantiles. The most common conditional models for discrete dependent variables are highly parametrized such as the Poisson model for count data or the ordered probit model for ordered data. These models have the advantages of being parsimonious in terms of parameters and easy to interpret. However, they impose strong homogeneity restrictions on the effects of the covariates. For instance, if a covariate increases the average outcome, then it must increase all the quantiles of the outcome distribution. Moreover, as pointed out by Winkelmann (2006), Poisson regression models imply a restrictive single crossing property on the sign of the estimated probability effects. These limitations are avoided by the distribution regression (DR) method (e.g., Williams and Grizzle (1972), Foresi and Peracchi (1995), and Chernozhukov et al. (2013)), which we employ in our paper. DR is a comprehensive tool for modeling and estimating the entire conditional distribution of any type of outcome (discrete, continuous, or mixed). DR allows the covariates to affect differently the outcome at different points of the distribution. The cost of this flexibility is that the DR parameters can be hard to interpret because they do not correspond to QEs. To overcome this problem, we report QEs computed as differences between the QFs of counterfactual distributions estimated by DR. These one-dimensional functions provide an intuitive summary of the effects of the (discrete or continuous) covariates. We therefore

propose to report DR-based estimates in conjunction with simultaneous confidence bands constructed using our projection method. We argue that this combination of our generic procedure with the DR model provides a comprehensive and practical approach for estimation QFs and QE functions with discrete data. (While we focus on DR in this paper, we emphasize that our projection method also combines well with classical parametric models such as Poisson regression and models alike, and so readers may find it useful for this reason alone).

Chernozhukov et al. (2013) also use DR to estimate counterfactual distributions but their results for QFs and QE functions apply only to continuous outcomes. They are based on asymptotic Gaussianity of the QFs and QE estimators. Their construction breaks down for discrete or mixed outcomes because the quantile operator is not smooth (Hadamard differentiable), which precludes the application of the delta method. Instead, we exploit the fact that the DR-based estimators of the DFs converge to Gaussian processes even for non-continuous outcomes. Then, we construct the confidence bands for the QFs and QE functions by inversion and Minkowski difference of the confidence bands for the DFs.

We apply our approach to two data sets, corresponding to two common types of discrete outcomes. In the first application, we exploit a large-scale randomized control trial in Oregon to estimate the distributional impact of insurance coverage on health care utilization measured by the number of doctor visits. Since this outcome is a count, we estimate the conditional DFs using both Poisson and distribution regressions. The Poisson regression clearly underestimates the probability of having zero visits as well as that of having a large number of visits. The more flexible DR finds a positive effect, especially at the upper tail of the distribution. This is an interesting empirical finding in its own right; it complements the mean regression analysis results reported in Finkelstein et al. (2012).

In the second application, we reanalyze the racial test score gap at the ages of eight

months and seven years. We find that while there is very little gap at eight months, a large gap arises at seven years. In addition, looking at the whole distribution, we uncover that the observed racial gap is widening in the upper tail of the distribution of test scores. The increase in the gap can be mostly explained by differences in observed characteristics between white and black children. These results complement and expand the findings of Fryer Jr and Levitt (2013) for the mean racial test score gap; our analysis is more complete, revealing what happens to the entire distribution.

The rest of the paper is organized as follows. Section 2 introduces our generic method to construct confidence bands for QFs and QE functions from simultaneous bands for DFs. Section 3 provides algorithms to construct simultaneous confidence bands for DFs based on bootstrap. Section 4 presents the two empirical applications.

2 Generic confidence bands

This section contains the main theoretical results of the paper. Our only assumption is the availability of simultaneous confidence bands for DFs. Since the seminal work of Kolmogoroff (1933), a variety of methods can be used to obtain these bands.¹ In Section 3 we describe specific algorithms that can be applied when the estimators of the DFs are known to be bootstrappable, which is often the case. In many cases the point and interval estimates of the DFs do not satisfy the logical monotonicity or range restrictions implied by the definition of a DF. In Section 2.1 we show how to impose these restrictions. In Section 2.2 we show how to invert confidence bands for DFs to obtain confidence bands for the corresponding QFs. In Section 2.3 we show how to combine the bands for two QFs to

¹The original Kolmogorov bands are actually conservative for discrete random variables, see Kolmogoroff (1941). Alternative methods, such as those described in Section 3, are asymptotically exact.

obtain a band for the difference between these QFs, the so-called QE function.

2.1 Confidence Bands for Distribution Functions

Let \mathcal{Y} be a closed subinterval of the extended real line. Let \mathbb{D} denote the set of weakly increasing functions, mapping \mathcal{Y} to $[0, 1]$. We will call the elements of this set “distribution functions”, albeit some of them need not be proper DFs. Let $y \mapsto F(y)$ in \mathbb{D} denote some target DF. This target could be a conditional DF, a marginal DF, or a counterfactual DF.

Definition 1 (Confidence Band of Level p). Given two functions $y \mapsto U(y)$ and $y \mapsto L(y)$ in the set \mathbb{D} such that $L \leq U$, pointwise, we define a band $I = [L, U]$ as the collection of intervals

$$I(y) = [L(y), U(y)], \quad y \in \mathcal{Y}.$$

We say that I covers F if $F \in I$ pointwise, namely $F(y) \in I(y)$ for all $y \in \mathcal{Y}$. If U and L are some data-dependent bands, we say that $I = [L, U]$ is a confidence band for F of level p , if I covers F with probability at least p . ■

In many applications the point estimates \hat{F} and interval estimates $[L', U']$ for the target distribution F do not satisfy logical monotonicity or range restrictions, namely they do not take values in the set \mathbb{D} . Given such an ordered triple $L' \leq \hat{F} \leq U'$, we can always transform it into another ordered triple $L \leq \check{F} \leq U$ that obeys the logical monotonicity and shape restrictions. For example, we can set

$$\check{F} = \mathcal{S}(\hat{F}), \quad L = \mathcal{S}(L'), \quad U = \mathcal{S}(U'), \quad (2.1)$$

where \mathcal{S} is the shaping operator that given a function $y \mapsto f(y)$ yields a mapping $y \mapsto \mathcal{S}(f)(y) \in \mathbb{D}$ with

$$\mathcal{S}(f) = \mathcal{M}(0 \vee f \wedge 1),$$

where the maximum and minimum are taken pointwise, and \mathcal{M} is the rearrangement operator that given a function $f : \mathcal{Y} \mapsto [0, 1]$ yields a map $y \mapsto \mathcal{M}(f)(y) \in \mathbb{D}$.²

The *rearrangement operator* is defined as follows. Let T be a countable subset of \mathcal{Y} . In our leading case where f is the distribution function of a discrete random variable Y , we can choose T as the support of Y and extend f to \mathcal{Y} by constant interpolation, yielding a step function as the distribution of Y on \mathcal{Y} . If f is a distribution function of a continuous or mixed random variable Y , we can set T as a grid of values covering the support of Y where we evaluate f and extend f to \mathcal{Y} by linear interpolation. Given a $f : T \mapsto [0, 1]$, we first consider $\mathcal{M}f$ as a vector of sorted values of the set $\{f(t) : t \in T\}$, where the sorting is done in a non-decreasing order. Since T is an ordered set of the same cardinality as $\mathcal{M}f$, we can assign the elements of $\mathcal{M}f$ to T in one-to-one manner: to the k -th smallest element of T we assign the k -th smallest element of $\mathcal{M}f$. The resulting mapping $t \mapsto \mathcal{M}f(t)$ is the rearrangement operator. We can extend the rearranged function $\mathcal{M}f$ to \mathcal{Y} by constant or linear interpolation as we described above.

The following lemma shows that shape restrictions *improve* the finite-sample properties of the estimators and confidence bands.

Lemma 1 (Shaping Improves Point and Interval Estimates). *The shaping operator \mathcal{S}*

(a) *is weakly contractive under the max distance:*

$$\|\mathcal{S}(A) - \mathcal{S}(B)\|_\infty \leq \|A - B\|_\infty, \quad \text{for any } A, B: T \rightarrow [0, 1],$$

(b) *is shape-neutral,*

$$\mathcal{S}(F) = F \text{ for any } F \in \mathbb{D},$$

²Other monotonization operators, such as the projection on the set of weakly increasing functions, can also be used, as we remark further below.

(c) and preserves the partial order:

$$A \leq B \implies \mathcal{S}(A) \leq \mathcal{S}(B), \quad \text{for any } A, B: T \rightarrow [0, 1].$$

Consequently,

1. the re-shaped point estimate constructed via (2.1) is weakly better than the initial estimate under the max distance:

$$\|\check{F} - F\|_\infty \leq \|\hat{F} - F\|_\infty,$$

2. the re-shaped confidence band constructed via (2.1) has weakly better coverage than the initial confidence band:

$$\mathbb{P}(L' \leq F \leq U') \leq \mathbb{P}(L \leq F \leq U),$$

3. and the re-shaped confidence band is weakly shorter than the original confidence bands under the max distance,

$$\|U - L\|_\infty \leq \|U' - L'\|_\infty.$$

Proof. The result follows from Chernozhukov et al. (2009). ■

The band $[L, U]$ is therefore weakly better than the original band $[L', U']$, in the sense that coverage is preserved while the width of the confidence band is weakly shorter.

Remark 1 (Isotonization is Another Option). An alternative to the rearrangement is the isotonization, which projects a given function on the set of weakly increasing functions that map T to $[0, 1]$. This also has the improving properties stated in Lemma 1. In fact any convex combination between isotonization and rearrangement has the improving properties stated in Lemma 1. ■

Remark 2 (Shape Restrictions on Confidence Bands by Intersection). An alternative way of imposing shape restrictions on the confidence band, is to intersect the initial band $[L', U']$ with WI, the set of weakly increasing functions that map T to $[0, 1]$. That is, we simply set

$$[L^I, U^I] = \text{WI} \cap [L', U'] = \{w \in \text{WI} : L'(y) \leq w(y) \leq U'(y), \quad \forall y \in T\}.$$

Thus, U^I is the greatest weakly increasing minorant of $0 \vee U' \wedge 1$ and L^I is the smallest weakly increasing majorant of $0 \vee L' \wedge 1$. This approach gives the tightest confidence bands, in particular

$$[L^I, U^I] \subseteq [L, U].$$

However, this construction might be less robust to misspecification than the rearrangement. For example, imagine that the target function F is not monotone, i.e. $F \notin \mathbb{D}$. This situation might arise when F is the probability limit of some estimator \hat{F} that is inconsistent for the DF due to misspecification. If the confidence band $[L', U']$ is sufficiently tight, then we can end up with an empty intersection band, $[L^I, U^I] = \emptyset$. By contrast $[L, U]$ is non-empty and covers the reshaped target function $F^* = \mathcal{S}(F) \in \mathbb{D}$. ■

2.2 Confidence Bands for Quantile Functions

Here we discuss the construction of confidence bands for the left-inverse function of F , $a \mapsto F^{\leftarrow}(a)$, which we can call the “quantile function” of F .

Definition 2 (Left Inverse). Given a function $y \mapsto G(y)$ in \mathbb{D} , we define its left inverse $a \mapsto G^{\leftarrow}(a)$ by:

$$G^{\leftarrow}(a) := \inf\{y \in \mathcal{Y} : G(y) \geq a\} \wedge \sup \mathcal{Y}, \quad a \in [0, 1],$$

where $\sup \mathcal{Y}$ is defined as $\sup \mathcal{Y} := \sup\{y \in \mathcal{Y}\}$. ■

The following theorem provides a confidence band I^\leftarrow for the QF F^\leftarrow based on a generic confidence band I for F .

Theorem 1 (Generic Bands for Quantile Functions). *Consider a distribution function F and band functions L and U in the class \mathbb{D} .*

1. *If F is covered by the band $I := [L, U]$, then the quantile function F^\leftarrow is covered by the band I^\leftarrow defined by*

$$I^\leftarrow(a) := [U^\leftarrow(a), L^\leftarrow(a)], \quad a \in [0, 1].$$

2. *Thus if the distribution function F is covered by I with probability p , then the quantile function F^\leftarrow is covered by I^\leftarrow with probability p .*

Proof. The result is immediate from the definition of the left inverse: For any $a \in [0, 1]$, since $L(y) \leq F(y)$ for each $y \in \mathcal{Y}$ and $F, L \in \mathbb{D}$,

$$\begin{aligned} F^\leftarrow(a) &= \inf\{y \in \mathcal{Y} : F(y) \geq a\} \wedge \sup \mathcal{Y} \\ &\leq \inf\{y \in \mathcal{Y} : L(y) \geq a\} \wedge \sup \mathcal{Y} = L^\leftarrow(a). \end{aligned}$$

Analogously, conclude that $F^\leftarrow(a) \geq U^\leftarrow(a)$. ■

This is the first main result of our paper. It shows that the band I^\leftarrow can literally be obtained by applying the left inverse transformation to the band I . We can narrow I^\leftarrow without affecting its coverage by exploiting the support restriction that the quantiles can only take the values of the underlying random variable. This is relevant when the variable of interest is discrete as in the applications presented in Section 4. Suppose that T is the support of the random variable with DF F . Then it makes sense to exploit the support restriction that $F^\leftarrow(a) \in T$ by intersecting the confidence bands for $F^\leftarrow(a)$ with T . Clearly, this will not affect the coverage properties of the bands.

Corollary 1 (Imposing Support Restrictions). *Consider the set \tilde{I}^{\leftarrow} defined by pointwise intersection of I^{\leftarrow} with T , namely $\tilde{I}^{\leftarrow}(a) := I^{\leftarrow}(a) \cap T$. Then, $\tilde{I}^{\leftarrow} \subseteq I^{\leftarrow}$ pointwise, and if I^{\leftarrow} covers F^{\leftarrow} then so does \tilde{I}^{\leftarrow} .*

The corollary is immediate because pointwise intersection of I^{\leftarrow} with the set T does not change the coverage property, since F^{\leftarrow} only takes values in T .

Figure 2 illustrates the construction of bands using Theorem 1 and Corollary 1. The left panel shows a DF $F : [0, 10] \mapsto [0, 1]$ covered by a band $I = [L, U]$. The middle panel shows that the inverse map $F^{\leftarrow} : [0, 1] \mapsto [0, 10]$ is covered by the inverted band $I^{\leftarrow} = [U^{\leftarrow}, L^{\leftarrow}]$. The band I^{\leftarrow} is easy to obtain by rotating and flipping I , but does not exploit the fact that the support of the variable with distribution F in this example is the set $T = \{0, 1, \dots, 10\}$. By intersecting I^{\leftarrow} with T for each $a \in [0, 1]$ we obtain in the right panel the band \tilde{I}^{\leftarrow} which reflects the support restrictions.

2.3 Generic Confidence Bands for Quantile Effects

The quantile effect (QE) function $a \mapsto \Delta_{j,m}(a)$ is the difference between the QFs of two random variables with DFs F_j and F_m and support sets T_j and T_m , i.e.

$$\Delta_{j,m}(a) := F_j^{\leftarrow}(a) - F_m^{\leftarrow}(a), \quad a \in [0, 1].$$

Our next goal is to construct simultaneous confidence bands that jointly cover the DFs $(F_k)_{k \in \mathcal{K}}$, where \mathcal{K} is a finite set, e.g., $\mathcal{K} = \{0, 1\}$, the corresponding QFs, and the QE functions $(\Delta_{jm})_{(j,m) \in \mathcal{K}^2}$.

Specifically, suppose we have the confidence bands $(I_k^{\leftarrow})_{k \in \mathcal{K}}$, which jointly cover the DFs $(F_k)_{k \in \mathcal{K}}$ with probability at least p . For example, we can construct these bands

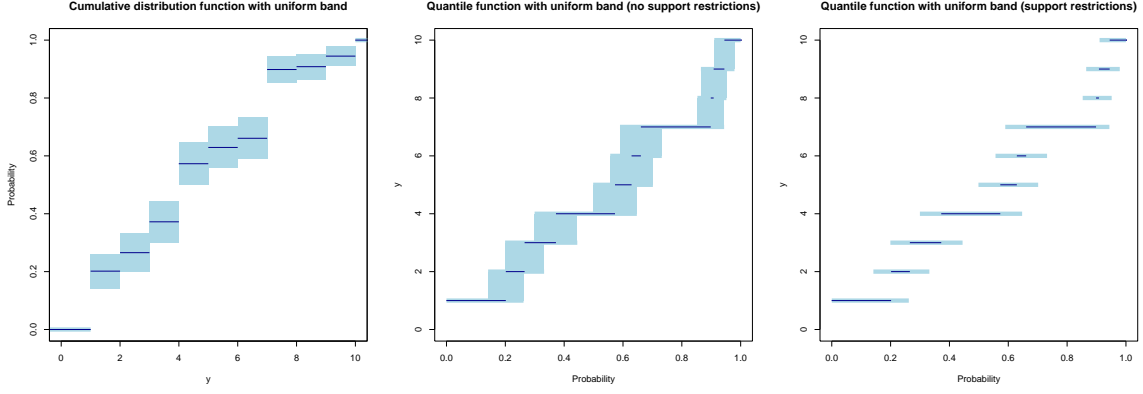


Figure 2: Construction of the bands for QF using Theorem 1 and Corollary 1. Left: the distribution function F (dark line) and confidence bands I (light rectangles). Middle: the quantile function F^{\leftarrow} and the confidence band I^{\leftarrow} . Right: The support-restricted confidence band \tilde{I}^{\leftarrow}

using Theorem 1 in conjunction with the Bonferroni inequality.³ The generic Algorithm 2 presented in Section 3 provides a construction of a joint confidence band that is not conservative. Then we can convert these bands to confidence bands for $\Delta_{j,m}$ for $(j, k) \in \mathcal{K}^2$ by taking the pointwise Minkowski difference \ominus of each of the pairs of the two bands, viewed as sets. Recall that the Minkowski difference between two subsets V and U of a vector space is $V \ominus U := \{v - u : v \in V, u \in U\}$. Note that if V and U are intervals $[v_1, v_2]$ and $[u_1, u_2]$, then

$$V \ominus U = [v_1, v_2] \ominus [u_1, u_2] = [v_1 - u_2, v_2 - u_1].$$

This greatly simplifies the practical computation of the bands.

³The joint coverage of two confidence bands with marginal coverage probabilities \tilde{p} is at least $p = 2\tilde{p} - 1$ by Bonferroni inequality.

Theorem 2 (Generic Bands for Quantile Effect Functions). *Consider the distribution functions $(F_k)_{k \in \mathcal{K}}$ and the band functions $(I_k := [L_k, U_k])_{k \in \mathcal{K}}$ in the class \mathbb{D} .*

1. *If F_j is covered by I_j and F_m is covered by I_m for $(j, m) \in \mathcal{K}^2$, then the quantile effect function $\Delta_{j,m} = F_j^{\leftarrow} - F_m^{\leftarrow}$ is covered by the band $I_{\Delta(j,m)}^{\leftarrow} = [U_j^{\leftarrow}, L_j^{\leftarrow}] - [U_m^{\leftarrow}, L_m^{\leftarrow}]$, where the minus operator is defined by a pointwise Minkowski difference:*

$$I_{\Delta(j,m)}^{\leftarrow}(a) := [U_j^{\leftarrow}(a), L_j^{\leftarrow}(a)] \ominus [U_m^{\leftarrow}(a), L_m^{\leftarrow}(a)], \quad a \in [0, 1].$$

2. *If the distribution functions $(F_k)_{k \in \mathcal{K}}$ are jointly covered by $(I_k)_{k \in \mathcal{K}}$ with probability p , then*

$$\mathbb{P}(F_k \in [L_k, U_k], F_k^{\leftarrow} \in I_k^{\leftarrow}, F_j^{\leftarrow} - F_m^{\leftarrow} \in I_{\Delta(j,m)}^{\leftarrow}; \text{ for all } (k, j, m) \in \mathcal{K}^3) = p.$$

Proof. Claim 1 is immediate from the definition of the Minkowski difference. Claim 2 follows because the event $\cap_{k \in \mathcal{K}} \{F_k \in [L_k, U_k]\}$ implies the event $\cap_{k \in \mathcal{K}} \{F_k^{\leftarrow} \in I_k^{\leftarrow}\}$ by Theorem 1, which implies the event $\cap_{(j,m) \in \mathcal{K}^2} \{F_j^{\leftarrow} - F_m^{\leftarrow} \in I_{\Delta(j,m)}^{\leftarrow}\}$. \blacksquare

This is the second main result of the paper. It shows that valid confidence bands for the QE function can be obtained by taking the Minkowski difference between the two confidence bands for the corresponding QFs. As in Theorem 1, we can narrow the band I_{Δ}^{\leftarrow} without affecting coverage by imposing support restrictions as demonstrated in Corollary 2.

Corollary 2 (Imposing Support Restrictions). *For $(j, m) \in \mathcal{K}^2$, consider the bands $\tilde{I}_{\Delta(j,m)}^{\leftarrow} = \tilde{I}_j^{\leftarrow} - \tilde{I}_m^{\leftarrow}$ defined by:*

$$\tilde{I}_{\Delta(j,m)}^{\leftarrow}(a) := \tilde{I}_j^{\leftarrow}(a) \ominus \tilde{I}_m^{\leftarrow}(a), \quad \tilde{I}_k^{\leftarrow}(a) := \{[U_k^{\leftarrow}(a), L_k^{\leftarrow}(a)] \cap T_k\}, \quad k \in \mathcal{K}.$$

Then $\tilde{I}_{\Delta(j,m)}^{\leftarrow} \subseteq I_{\Delta(j,m)}^{\leftarrow}$, and if $I_{\Delta(j,m)}^{\leftarrow}$ covers $\Delta_{j,m}$ then so does $\tilde{I}_{\Delta(j,m)}^{\leftarrow}$.

Remark 3 (Joint Support Restrictions). The band $\tilde{I}_{\Delta(j,m)}^{\leftarrow}$ can be further narrowed if the two random variables with distributions F_j and F_m have restrictions in their joint support T_{jm} , i.e., $T_{jm} \neq T_j \times T_m = \{(t_j, t_m) : t_j \in T_j, t_m \in T_m\}$. In this case we can drop all the elements d from $\tilde{I}_{\Delta(j,m)}^{\leftarrow}$ that cannot be formed as $d = t_j - t_m$ for some $(t_j, t_m) \in T_{jm}$. For example, let $T_j = T_m = \tilde{I}_{\Delta(j,m)}^{\leftarrow} = \{0, 1, 2\}$, then we can drop 2 from $\tilde{I}_{\Delta(j,m)}^{\leftarrow}$ if $(0, 2) \notin T_{jm}$. ■

Remark 4 (Similarity). Theorem 2 shows that our generic method of constructing bands carries over the similarity (non-conservativeness) of the bands for the DFs to the simultaneous bands for the QFs and QE functions. Moreover, our construction is optimal in the sense that if we want to simultaneously cover all the distribution, quantile and QE functions of interest, it is not possible to construct uniformly shorter bands while preserving the joint coverage rate once all the joint support restrictions are imposed.

It is common to report at the same time several quantile and QE functions. For instance, Figures 5 and 6 provide three different QFs (two observed and one counterfactual) and the differences between these functions, which are all of interest. Theorem 2 (together with Corollary 4 for the asymptotic similarity of the bands for the DFs) shows that our bands jointly cover asymptotically all these functions with probability p . This allows for a transparent and honest assessment of hypotheses about these functions.

On the other hand, when the goal is to cover only a single QE function independently from the other functions, then our band for this function can be conservative. This is due to the projection implicit in the application of the Minkowski difference and is the price to pay for the joint uniform coverage property. However, our empirical results clearly demonstrate the usefulness of these bands that allow for testing hypotheses that could not be considered using existing methods. We are not aware of any generic method to construct nonconservative bands for QE functions of discrete outcomes. ■

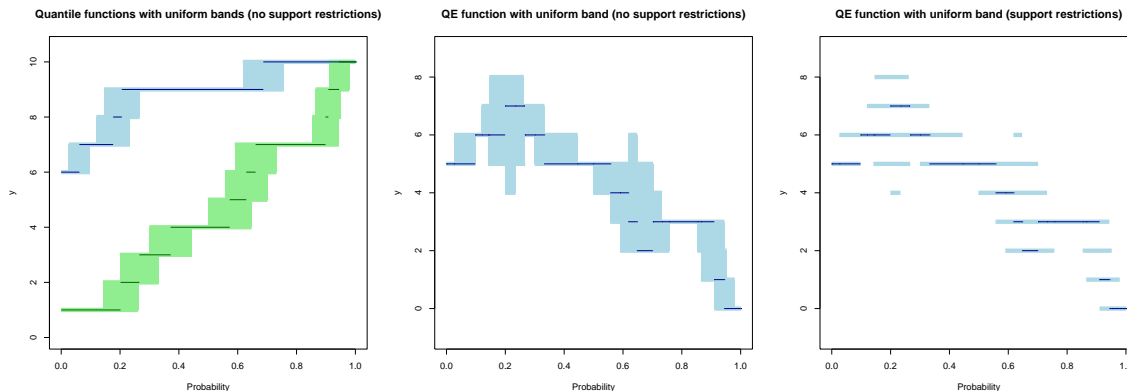


Figure 3: Construction of the bands for QE functions using Theorem 2 and Corollary 2. Left: quantile functions F_0^{\leftarrow} and F_1^{\leftarrow} and confidence bands I_0^{\leftarrow} and I_1^{\leftarrow} . Middle: the QE function Δ and the confidence band \bar{I}_Δ without support restrictions. Right: the QE function Δ and the confidence band \tilde{I}_Δ with support restrictions.

Figure 3 illustrates the construction of the bands for QE functions using Theorem 2 and Corollary 2. The left panel shows the bands I_0^{\leftarrow} and I_1^{\leftarrow} for the QFs F_0^{\leftarrow} and F_1^{\leftarrow} . The middle panel shows the band $I_{\Delta(1,0)}$ for the QE function $\Delta_{1,0} = F_1^{\leftarrow} - F_0^{\leftarrow}$, obtained by taking the Minkowski difference of I_1^{\leftarrow} and I_0^{\leftarrow} . The right panel shows the confidence band $\tilde{I}_{\Delta(1,0)}$ for the QE function $\Delta_{1,0}$ resulting from imposing the support restrictions. As the Theorem 2 proves, the QE function $\Delta_{1,0}$ is covered by the band $I_{\Delta(1,0)}$.

3 Bootstrap Algorithms for Constructing Simultaneous Confidence Bands

In Section 2 we assumed the existence of simultaneous confidence bands for DFs. Here we describe precise algorithms that are shown to provide asymptotically valid simultaneous bands for any bootstrappable estimator of the DF. Many commonly used estimators of the DF are bootstrappable under suitable conditions. For example, Theorems 4.2 and 5.2 of Chernozhukov et al. (2013) give conditions for bootstrap consistency for the DR based estimators that we use in the empirical applications. Maximum likelihood estimators, such as the Poisson regression that we use as a benchmark in the first application, are also bootstrappable under weak differentiability conditions, see Arcones and Giné (1992). We note that if the data are discrete, these existing results yield valid uniform bands for the DFs but *cannot* be used to construct uniform bands for the QFs and QE functions.

Algorithm 1 provides uniform confidence bands that asymptotically cover the DF F and the QF F^{\leftarrow} with probability p :

Algorithm 1 (Bootstrap Algorithm for Confidence Bands for F and F^{\leftarrow}).

1. Obtain many bootstrap draws of the estimator \hat{F} ,

$$\hat{F}^{*(j)}, \quad j = 1, \dots, B$$

where the index j enumerates the bootstrap draws and B is the number of bootstrap draws (e.g., $B = 1,000$).

2. For each y in T , compute the robust standard error of $\hat{F}(y)$,

$$\hat{s}(y) = (\hat{Q}(.75, y) - \hat{Q}(.25, y)) / (\Phi^{\leftarrow}(.75) - \Phi^{\leftarrow}(.25)),$$

where $\hat{Q}(\alpha, y)$ denotes the empirical α -quantile of the bootstrap sample $(\hat{F}^{*(j)}(y))_{j=1}^B$, and Φ^{\leftarrow} denotes the inverse of the standard normal distribution.

3. Compute the critical value

$$c(p) = p\text{-quantile of } \left\{ \max_{y \in T} |\hat{F}(y)^{*(j)} - \hat{F}(y)| / \hat{s}(y) \right\}_{j=1}^B.$$

4. Construct a preliminary uniform confidence band $[L', U']$ for F of level p via: $[L'(y), U'(y)] = [\hat{F}(y) \pm c(p)\hat{s}(y)]$ for each $y \in T$. Impose the shape restrictions on \hat{F} , L' and U' by setting:

$$\check{F} = \mathcal{S}(\hat{F}), \quad [L, U] = [\mathcal{S}(L'), \mathcal{S}(U')].$$

Report $I = [L, U]$ as a p -level uniform confidence band for F .

5. Report the inverted band $I^{\leftarrow} = [U^{\leftarrow}, L^{\leftarrow}]$ or support restricted inverted band $\tilde{I}^{\leftarrow} = I^{\leftarrow} \cap T$ as a p -level uniform confidence band for F^{\leftarrow}

In the first step we bootstrap the estimator of the DF of interest. In the second step we estimate the pointwise standard errors using the bootstrap. In the third step, we compute the p -quantile of the bootstrap draws of the weighted centered Kolmogorov–Smirnov statistic, and construct preliminary and shape-restricted uniform bands. Lemma 1 shows that the shape restrictions bring about finite-sample improvements. In the last step we invert the band for the DF to obtain a uniform band for the QF, as justified by Theorem 1; and, if needed, we can impose the support conditions.

Remark 5 (Step (1)). There are multiple ways to obtain the bootstrap draws of \hat{F} . A

generic resampling procedure is the exchangeable bootstrap (Praestgaard and Wellner, 1993; van der Vaart and Wellner, 1996), which recomputes \hat{F} using sampling weights drawn independently from the data. This procedure incorporates many popular bootstrap schemes as special cases by a suitable choice of the distribution of the weights. For example, the empirical bootstrap corresponds to multinomial weights, and the weighted or Bayesian bootstrap corresponds to standard exponential weights. Exchangeable bootstrap can also accommodate dependences or clustering in the data by drawing the same weight for all the observations that belong to the same cluster (Sherman and Cessie, 1997; Cheng et al., 2013). For example, in the application of Section 4.3 we draw the same weights for all the individuals of the same household. ■

The following result provides a theoretical justification for this algorithm. In what follows, let $\ell^\infty(\mathcal{Y})$ denote the metric space of bounded function from \mathcal{Y} to \mathbb{R} .

Corollary 3 (Validity of Algorithm 1). *Suppose that the rescaled DF estimator $a_n(\hat{F} - F)$ converges in law in $\ell^\infty(\mathcal{Y})$ to a Gaussian process G , having zero mean and a non-degenerate variance function, for some sequence of constants $a_n \rightarrow \infty$ as $n \rightarrow \infty$, where n is some index (typically the sample size). Suppose that a bootstrap method can consistently approximate the limit law of $a_n(\hat{F} - F)$, namely the distance between the law of $a_n(\hat{F}^* - \hat{F})$ conditional on data, and that of G , converges to zero in probability as $n \rightarrow \infty$. The distance is the bounded Lipschitz metric that metrizes weak convergence. Then,*

$$\lim_{n \rightarrow \infty} \mathbb{P}(F \in I, F^{\leftarrow} \in \tilde{I}^{\leftarrow}) = p.$$

Proof. Lemma SA.1 of Chernozhukov et al. (2013) implies that $\lim_{n \rightarrow \infty} \mathbb{P}(F \in [L', U']) = p$. The result then follows from Lemma 1, Theorem 1 and Corollary 1. ■

Algorithm 2 provides simultaneous confidence bands that asymptotically jointly cover

the DFs $(F_k)_{k \in \mathcal{K}}$, where \mathcal{K} is a finite set, e.g., $\mathcal{K} = \{0, 1\}$, the corresponding QFs, and the QE functions $F_j^{\leftarrow} - F_k^{\leftarrow}$ for all $(j, k) \in \mathcal{K}^2$, with probability p .

Algorithm 2 (Bootstrap Algorithm for Confidence Bands for QFs and QEs).

1. Obtain many bootstrap draws of the estimator $(\hat{F}_k)_{k \in \mathcal{K}}$,

$$(\hat{F}_k^{*(j)})_{k \in \mathcal{K}}, \quad j = 1, \dots, B,$$

where the index j enumerates the bootstrap draws and B is the number of bootstrap draws (e.g., $B = 1,000$).

2. For each y in T and k in \mathcal{K} , compute the robust standard error of $\hat{F}_k(y)$:

$$\hat{s}_k(y) = (\hat{Q}_k(.75, y) - \hat{Q}_k(.25, y)) / (\Phi^{\leftarrow}(.75) - \Phi^{\leftarrow}(.25)),$$

where $\hat{Q}_k(\alpha, y)$ denotes the empirical α -quantile of the bootstrap sample $(\hat{F}_k^{*(j)}(y))_{j=1}^B$, and Φ^{\leftarrow} denotes the inverse of the standard normal distribution.

3. Compute the critical value

$$c(p) = p\text{-quantile of } \left\{ \max_{y \in T, k \in \mathcal{K}} |\hat{F}_k^{*(j)}(y) - \hat{F}_k(y)| / \hat{s}_k(y) \right\}_{j=1}^B.$$

4. Construct preliminary joint confidence bands $([L'_k, U'_k])_{k \in \mathcal{K}}$ for $(F_k)_{k \in \mathcal{K}}$ of level p as

$$[L'_k(y), U'_k(y)] = [\hat{F}_k(y) \pm c(p)\hat{s}_k(y)], \quad y \in T, \quad k \in \mathcal{K}.$$

For each $k \in \mathcal{K}$ impose the shape restrictions on \hat{F}_k , L'_k and U'_k by setting:

$$\check{F}_k = \mathcal{S}(\hat{F}_k), \quad [L_k, U_k] = [\mathcal{S}(L'_k), \mathcal{S}(U'_k)].$$

5. Report $(I_k)_{k \in \mathcal{K}} = ([L_k, U_k])_{k \in \mathcal{K}}$ as p -level simultaneous confidence bands for $(F_k)_{k \in \mathcal{K}}$. Report $(I_k^{\leftarrow})_{k \in \mathcal{K}} = ([U_k^{\leftarrow}, L_k^{\leftarrow}])_{k \in \mathcal{K}}$ or the support-restricted version $(\tilde{I}_k^{\leftarrow})_{k \in \mathcal{K}} = (I_k^{\leftarrow} \cap T_k)_{k \in \mathcal{K}}$ as p -level simultaneous confidence bands for $(F_k^{\leftarrow})_{k \in \mathcal{K}}$.
6. Report $I_{\Delta(j,k)}^{\leftarrow} = I_j^{\leftarrow} \ominus I_k^{\leftarrow}$ or the support-restricted version $\tilde{I}_{\Delta(j,k)}^{\leftarrow} = \tilde{I}_j^{\leftarrow} \ominus \tilde{I}_k^{\leftarrow}$ as p -level simultaneous confidence bands for the quantile effect functions $F_j^{\leftarrow} - F_k^{\leftarrow}$ for all $(j, k) \in \mathcal{K}^2$.

Compared to the first algorithm the crucial differences are the following. First, we bootstrap jointly all the estimators of the DFs in step (1). In our applications these estimators are not independent such that it is important to obtain jointly the bootstrap draws of all them. Secondly, the Kolmogorov-Smirnov maximal t -statistic in step (3) is computed over all distributions F_k with $k \in \mathcal{K}$. This ensure that the confidence band in step (4) covers jointly all the DFs with probability p . In the last step we obtain the confidence band for the QE functions by taking Minkowski differences, as justified by Theorem 2.

Let $|\mathcal{K}|$ denote the cardinality of the set \mathcal{K} .

Corollary 4 (Validity of Algorithm 2). *Suppose that the rescaled DF estimators $\{a_n(\hat{F}_k - F_k)\}_{k \in \mathcal{K}}$ converge in law in $\ell^\infty(\mathcal{Y})^{|\mathcal{K}|}$ to a Gaussian process $(G_k)_{k \in \mathcal{K}}$, having zero mean and a non-degenerate variance function, for some sequence of constants $a_n \rightarrow \infty$ as $n \rightarrow \infty$, where n is some index (typically the sample size). Suppose that a bootstrap method can consistently approximate the limit law of $\{a_n(\hat{F}_k - F_k)\}_{k \in \mathcal{K}}$, namely the distance between the law of $\{a_n(\hat{F}_k^* - \hat{F}_k)\}_{k \in \mathcal{K}}$ conditional on data, and that of $(G_k)_{k \in \mathcal{K}}$, converges to zero in probability as $n \rightarrow \infty$. The distance is the bounded Lipschitz metric that metrizes weak*

convergence. Then, the confidence bands constructed by Algorithm 2 have the following covering property:

$$\lim_{n \rightarrow \infty} \mathbb{P}(F_k \in I_k, F_k^{\leftarrow} \in \tilde{I}_k^{\leftarrow}, F_j^{\leftarrow} - F_m^{\leftarrow} \in \tilde{I}_{\Delta(j,m)}^{\leftarrow}; \text{ for all } (k, j, m) \in \mathcal{K}^3) = p.$$

Proof. Lemma SA.1 of Chernozhukov et al. (2013) implies that $\lim_{n \rightarrow \infty} \mathbb{P}(\cap_{k \in \mathcal{K}} \{F_k \in [L'_k, U'_k]\}) = p$. The result then follows from Corollary 3, Theorem 2 and Corollary 2. ■

4 Applications to Distribution Regression Analysis of Discrete Data

In this section we apply our approach to two data sets, corresponding to two common types of discrete outcomes. In both cases we use the distribution regression model and obtain QE as differences between counterfactual distributions. For this reason, we first introduce the specific methods and then present both empirical illustrations.

4.1 Distribution Regression

In the absence of covariates, the empirical DF is a minimal sufficient statistic for a non-parametric marginal DF. Distribution regression (DR) generalizes this concept to a conditional DF like OLS generalizes the univariate mean to the conditional mean function. The key, simple observation underlying DR is that the conditional distribution of the outcome Y given the covariates X can be expressed as $F_{Y|X}(y | x) = \mathbb{E}[1\{Y \leq y\} | X = x]$. Accordingly, we can always construct a collection of binary response variables, which record the events that the outcome Y falls bellow a set of thresholds T , i.e.,

$$1\{Y \leq y\}, \quad y \in T,$$

and use a binary regression model for each variable in this collection. This yields the DR model:

$$F_{Y|X}(y | x) = P(Y \leq y | X = x) = \Lambda_y(B(x)' \beta(y)), \quad (4.1)$$

where $\Lambda_y(\cdot)$ is a known link function which is allowed to change with the threshold level y ; $B(x)$ is a vector of transformations of x with good approximating properties such as polynomials, B-splines, and interactions; and $\beta(y)$ is an unknown vector of parameters. Knowledge of the function $y \mapsto \beta(y)$ implies knowledge of the distribution of Y conditional on X . The DR model is flexible in the sense that, for any given link function, we can approximate the conditional DF arbitrarily well by using a rich enough set of transformations of the original covariates $B(x)$. In the extreme case when X is discrete and $B(x)$ is fully saturated, the estimated conditional distribution is numerically equal to the empirical DF in each cell of X for any monotonic link function. When $B(x)$ is not fully saturated, one can choose a DF such as the normal or logistic as the link function to guarantee that the model probabilities lie between 0 and 1.

DR nests a variety of classical models such as the normal regression, the Cox proportional hazard, ordered logit, ordered probit, Poisson regression, as well as other generalized linear models. Example 1 shows the inclusion of the Poisson regression model which we use as a benchmark in our first empirical application. In what follows we set $B(x) = x$ to lighten the notation without loss of generality.

Example 1. Let Y be a nonnegative integer-valued outcome and X a vector of covariates. The Poisson regression model assumes that the probability mass function of Y conditional on X is

$$f_{Y|X}(y | x) = \frac{\exp(x'\beta)^y \exp(-\exp(x'\beta))}{y!} \text{ for } y = \{0, 1, 2, \dots\}.$$

The corresponding conditional distribution is:

$$F_{Y|X}(y | x) = \sum_{k=0}^y \frac{\exp(x'\gamma)^k \exp(-\exp(x'\beta))}{k!} = Q(y, \exp(x'\beta)),$$

where Q is the incomplete gamma function. Thus, the Poisson regression can be seen as a special case of a DR model with exponentiated incomplete gamma link function,

$$\Lambda_y(u) = Q(y, \exp u), \quad (4.2)$$

and parameter function $y \mapsto \beta(y)$ that does not vary with y , i.e. $\beta(y) = \beta$. The Poisson regression model therefore imposes strong homogeneity restrictions on the effect of the covariates at different parts of the distribution that are often rejected by the data (see, e.g., Section 4.3). ■

Assume that we have a sample $\{(Y_i, X_i) : i = 1, \dots, n\}$ of (Y, X) . The DR estimator of the conditional distribution is

$$\hat{F}_{Y|X}(y | x) = \Lambda_y(x' \hat{\beta}(y)), \quad y \in T,$$

where

$$\hat{\beta}(y) = \arg \max_{b \in \mathbb{R}^{\dim(X)}} \sum_{i=1}^n 1\{Y_i \leq y\} \ln [\Lambda_y(X_i' b)] + 1\{Y_i > y\} \ln [1 - \Lambda_y(X_i' b)].$$

Williams and Grizzle (1972) introduced DR in the context of ordered outcomes. Foresi and Peracchi (1995) applied this method to estimate the conditional distribution of excess return evaluated at a finite number of points. Chernozhukov et al. (2013) extended Williams and Grizzle (1972)'s definition to arbitrary outcomes and established functional central limit theorems and bootstrap validity results for DR as an estimator of the whole conditional distribution. One of the main advantages of DR is that it not only accommodates continuous but also discrete and mixed discrete continuous outcomes very naturally.

4.2 Marginal and Counterfactual Distributions

In the two applications that we present below there are two groups: the treated and control units in the first application, and the black and white children in the second application. We use DR to model and estimate the conditional distribution of the outcome in each group at each value of the covariates, that we denote by $F_{Y_0|X_0}(y | x)$ and $F_{Y_1|X_1}(y | x)$. The difference between these two high-dimensional DFs is, however, difficult to convey. Instead, we integrate these conditional distributions with respect to observed covariate distributions and compare the resulting marginal distributions.

For instance, in the first application, the marginal distribution

$$F_{\langle k \rangle}(y) := \int F_{Y_k|X_k}(y | x) dF_X(x),$$

where F_X is the distribution of X in the entire population including the treated and control units, represents the distribution of a potential outcome. When $k = 1$ is the outcome distribution that would be observed if every units were treated, and when $k = 0$ is the outcome distribution if every units were not treated. These two distributions are called counterfactual, since they do not arise as distributions from any observable population. They have nevertheless a causal interpretation as distributions of potential outcomes when the treatment is randomized conditionally on the control variables X .

Let $\hat{F}_{Y_k|X_k}$ denote the DR estimator of $F_{Y_k|X_k}$, $k \in \{0, 1\}$. We estimate $F_{\langle k \rangle}$ by the plugging-in rule, namely integrating $\hat{F}_{Y_k|X_k}$ with respect to the empirical distribution of X for treated and control units. For $k \in \{0, 1\}$,

$$\hat{F}_{\langle k \rangle}(y) := \frac{1}{n} \sum_{i=1}^n \hat{F}_{Y_k|X_k}(y | X_i).$$

We then report the empirical QE function:

$$\hat{\Delta}(a) := \hat{F}_{\langle 1 \rangle}^{\leftarrow}(a) - \hat{F}_{\langle 0 \rangle}^{\leftarrow}(a), a \in [0, 1].$$

Chernozhukov et al. (2013) derived joint functional central limit theorems for $(\hat{F}_{\langle 0 \rangle}, \hat{F}_{\langle 1 \rangle})$ and established bootstrap validity.⁴ We can thus use the algorithms in Section 3 to construct asymptotically valid simultaneous confidence bands for the counterfactual QFs $(F_{\langle 1 \rangle}^{\leftarrow}, F_{\langle 0 \rangle}^{\leftarrow})$ and the QE function $\Delta = F_{\langle 1 \rangle}^{\leftarrow} - F_{\langle 0 \rangle}^{\leftarrow}$.

Remark 6 (Continuous covariates). The proposed approach can also be used to analyze the effect of continuous covariates. For instance, we can compare the status quo QF with the QF that we would observe if everyone received Δd additional units of the continuous covariate of interest D , e.g. $\Delta d = 1$ for a unitary increase. Formally, assume that we are interested in the effect of a continuous variable D on the outcome Y while controlling for a vector of covariates X . We can define the counterfactual distribution

$$F_{\langle \Delta d \rangle}(y) := \int F_{Y|X,Z}(y | d + \Delta d, x) dF_{D,X}(d, x)$$

and the QE function $F_{\langle \Delta d \rangle}^{\leftarrow}(a) - F_{\langle 0 \rangle}^{\leftarrow}(a)$, where $F_{\langle 0 \rangle}$ is the marginal (status quo) distribution of Y . This experiment can be interpreted as an unconditional quantile regression. Also in this case, our methods provide valid confidence bands for the counterfactual quantile and QE functions. ■

4.3 Insurance coverage and health care utilization

Our first application illustrates the construction of the confidence bands using data from the Oregon health insurance experiment. In 2008, the state of Oregon initiated a limited expansion of its Medicaid program for uninsured low-income adults by offering insurance coverage to the lottery winners from a waiting list of 90,000 people (see www.nber.org/oregon for

⁴We refer to Chernozhukov et al. (2013) for a more comprehensive discussion of counterfactual distributions, their interpretation and the statistical properties of their estimators.

details). This experiment constitutes a unique opportunity to study the impact of insurance by means of a large-scale randomized controlled trial (e.g., Finkelstein et al., 2012; Baicker et al., 2013, 2014; Taubman et al., 2014).

We investigate the impact of insurance coverage on health care utilization as analyzed in Finkelstein et al. (2012, Section V) using a publicly available dataset.⁵ Detailed information about the dataset and descriptive statistics are available in Finkelstein et al. (2012) and the corresponding online appendix. We focus on one count outcome Y : the number of outpatient visits in the last six months, which was elicited via a large mail survey. After excluding individuals with missing information in any of the variables used in the analysis, the resulting sample consists of 23,441 observations. The top histogram in Figure 1 illustrates the discrete nature of our dependent variable. Almost 40% of the outcomes are zeros, more than 90% of the mass is concentrated between zero and five, but a few people have a greater number of visits.

Finkelstein et al. (2012) find a positive effect of winning the lottery on the number of outpatient visits.⁶ Their results are based on ordinary least squares (OLS) regressions, where the covariates X include household size, indicators for the survey wave, and interactions of the household size indicators and the survey wave. Although individuals were chosen randomly, these covariates are included as controls because the entire household for any selected individual became eligible to apply for insurance and the fraction of treated individuals varies across survey waves. We complement their findings by looking at the whole distribution of the number outpatient visits. We first estimate the conditional outcome distributions separately for the lottery winners and losers via Poisson regression and

⁵The data are available via: <http://www.nber.org/oregon/4.data.html>

⁶They label these effects intention-to-treat (ITT) effects and also report local average treatment effects (LATE) estimated using IV regressions. In this section, we focus on ITT effects.

DR. For DR, we use the exponentiated incomplete gamma link in (4.2) such that DR nests the Poisson regression as an exact special case. As explained in Section 4.2, we integrate the conditional outcome distributions with respect to the covariate distribution for both lottery winners and losers to obtain estimates of the counterfactual distributions $F_{\langle 1 \rangle}$ and $F_{\langle 0 \rangle}$.

The top panels of Figure 4 displays the DFs $\hat{F}_{\langle 1 \rangle}$ and $\hat{F}_{\langle 0 \rangle}$ estimated by the Poisson regression and DR. The corresponding QFs $\hat{F}_{\langle 1 \rangle}^{\leftarrow}$ and $F_{\langle 0 \rangle}^{\leftarrow}$ are displayed in both middle panels. Finally, the estimated QE functions, $\hat{F}_{\langle 1 \rangle}^{\leftarrow} - \hat{F}_{\langle 0 \rangle}^{\leftarrow}$, are plotted in the bottom panels. In all cases, the figure also shows 95% simultaneous confidence bands, constructed using Algorithm 1 and 2 with $B = 1,000$ Bayesian bootstrap draws that take into account the possible clustering of the observations at the household level. Reflecting the discrete nature of our outcome variables, we impose the support restrictions $T_0 = T_1 = \{0, 1, \dots\}$.

A comparison between the Poisson and DR results reveals striking differences. The Poisson model predicts a much lower mass at zero and a much thinner upper tail of the distribution for both groups. Indeed, these differences are statistically significant as the Poisson and DR simultaneous confidence bands for the DFs and QFs do not overlap for a large part of the support. A formal test rejects the equality of these distributions with a p-value below 0.001. Since the DR model with exponentiated incomplete gamma link nests the Poisson model, we conclude that the Poisson model is rejected by the data. For this reason, we focus the discussion on the DR results.

The simultaneous band for the QE function do not fully cover the zero-line such that we can reject the null hypothesis that winning the lottery has no effect on the number of outpatient visits. We can also reject the hypothesis that $F_{\langle 0 \rangle}$ first-order stochastically dominates $F_{\langle 1 \rangle}$ because the band for $F_{\langle 0 \rangle}^{\leftarrow}$ is strictly below the band for $F_{\langle 1 \rangle}^{\leftarrow}$ at some quantile indexes, but we cannot reject the opposite hypothesis. In other words, at no quantile index

the confidence band contains strictly negative effects while at some quantile indexes it contains strictly positive effects.

Health economists distinguish between the treatment effect on the extensive (whether to see a doctor) and intensive (the number of visits given at least one) margins. The first effect is easy to estimate: the probability of not seeing a doctor decreased significantly from 43% to 37% with the treatment. The effect on the intensive margin is more difficult to gauge because we do not observe both potential outcomes for any individual. If we assume that the individuals induced to see a doctor by the insurance coverage are not seriously sick and visit the doctor only once, then the effect on the intensive margin can also be seen in Figure 4: the effect from 0 to 1 visit represents the effect on the extensive margin and the effect on the rest of the distribution represent the effect on the intensive margin. Both effects are statistically significant. We note in particular that the quantile differences do not vanish at the top of the distribution.

The assumption made to justify this interpretation may be too strong and lead to an overestimation of the effect on the intensive margin. For instance, the doctor may find a serious problem and schedule other visits. Following Zhang and Rubin (2003) and Angrist et al. (2006), we can bound the effect on the intensive margin from below by assuming that patients who see a doctor anyway visit their doctor at least as often as patients who see a doctor only if insured. Under this weaker assumption, the effect on the intensive margin is bounded from below by the QE function obtained by keeping only observations with at least one visit. We also find a positive treatment effect with this method, which re-inforce the evidence of a positive effect among the existing users.

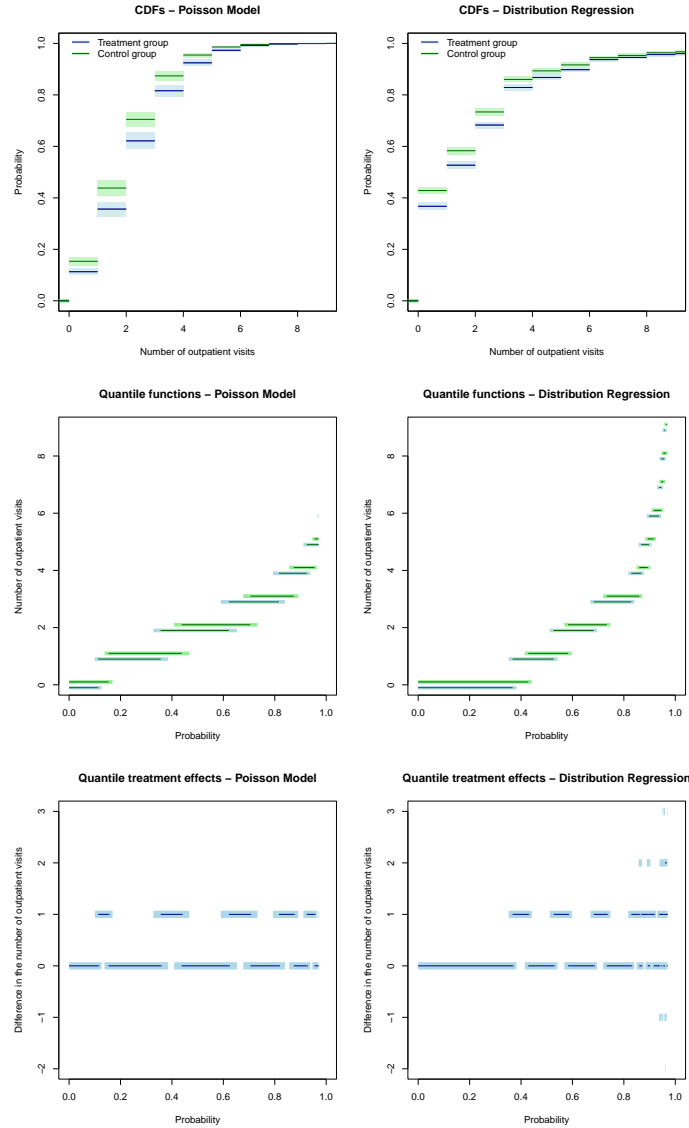


Figure 4: Effect of insurance coverage on the number of outpatient visits. DFs, quantile functions, and QTE estimated by Poisson regression and DR including support restricted 95% confidence bands. The lines of the quantile functions for the control group are slightly shifted upward to avoid overlapping with the quantile function for the treatment group.

4.4 Racial differences in mental ability of young children

As a second application, we reanalyze the racial IQ test score gap examined in Fryer and Levitt (2013). We use data from the US Collaborative Perinatal Project (CPP). These data contain information on children from 30,002 women who gave birth in 12 medical centers between 1959 and 1965. Our main outcomes of interest are the standardized test scores at the ages of eight months (Bayley Scale of Infant Development) and seven years (both Stanford-Binet and Wechsler Intelligence Test). In addition to the test score measures, the dataset contains a rich set of background characteristics for the children, X , including information on age, gender, region, socioeconomic status, home environment, prenatal conditions, and interviewer fixed effects. Fryer and Levitt (2013) provide a comprehensive description of the dataset and extensive descriptive statistics.

A key feature of the test scores is the discrete nature of their distribution. We observe only 76 and 128 different values for the standardized test scores at the ages of eight months and seven years, respectively. The middle and bottom panels of Figure 1 present the corresponding histograms. Note that each bar corresponds to exactly one value. For instance, at eight months, almost 12% of the observations have exactly the same score and 60% of the observations have one of the most frequent six values. This is a common feature of test scores, which are necessarily discrete because they are based on a finite number of questions.

To gain a better understanding of the causes of the observed black-white test score gap, we provide a distributional decomposition into explained and unexplained parts by observable background characteristics. Let $F_{\langle W|W \rangle}$ and $F_{\langle B|B \rangle}$ represent the observed test score DFs for white and black children, and $F_{\langle W|B \rangle}$ represent the counterfactual DF of test scores that would have prevailed for white children had they had the distribution of

background characteristics of black children, F_{X_B} , namely,

$$F_{\langle W|B \rangle}(y) := \int F_{Y_W|X_W}(y | x) dF_{X_B}(x). \quad (4.3)$$

With this counterfactual test score distribution it is possible to decompose the observed black-white test score gap into

$$F_{\langle W|W \rangle}^{\leftarrow} - F_{\langle B|B \rangle}^{\leftarrow} = [F_{\langle W|W \rangle}^{\leftarrow} - F_{\langle W|B \rangle}^{\leftarrow}] + [F_{\langle W|B \rangle}^{\leftarrow} - F_{\langle B|B \rangle}^{\leftarrow}]. \quad (4.4)$$

where the first term in brackets corresponds is the composition effect due to differences in observable background characteristics and the second term is the unexplained difference.

We estimate $F_{\langle W|W \rangle}$ and $F_{\langle B|B \rangle}$ by the empirical test score distributions for white and black children, respectively. We estimate the counterfactual distribution $F_{\langle W|B \rangle}$ by the sample analog of (4.3) replacing $F_{Y_W|X_W}$ by the DR estimator for white children, and F_{X_B} by the empirical distribution of X for black children. We use the logistic link function for the DR, but the results using the linear link function or the normal link function are similar.

Figures 5 and 6 show the results for the eight months and seven years outcomes. The first panels show the observed and counterfactual QFs, $F_{\langle W|W \rangle}^{\leftarrow}$, $F_{\langle B|B \rangle}^{\leftarrow}$ and $F_{\langle W|B \rangle}^{\leftarrow}$. The second panels show the difference between the observed QFs, $F_{\langle W|W \rangle}^{\leftarrow} - F_{\langle B|B \rangle}^{\leftarrow}$. The third and fourth panels decompose these observed differences into the composition effect ($F_{\langle W|W \rangle}^{\leftarrow} - F_{\langle W|B \rangle}^{\leftarrow}$) and the unexplained component ($F_{\langle W|B \rangle}^{\leftarrow} - F_{\langle B|B \rangle}^{\leftarrow}$). The point estimates are shown with their respective 95% simultaneous confidence bands constructed using Algorithm 1 and 2 with $B = 1,000$ Bayesian bootstrap draws. The bands impose the restrictions that the supports of the test scores correspond to the observed values in the sample.

For eight months old children, we find very small differences between the test score distributions of black and white children. The black-white gap is positive at the lower

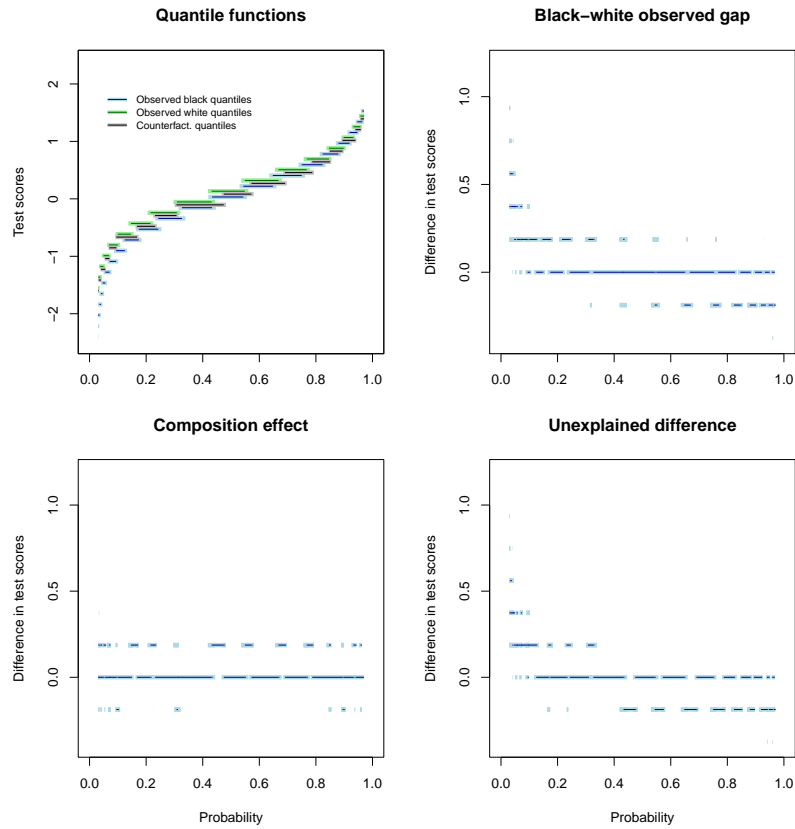


Figure 5: Decomposition of observed racial differences in mental ability of young children; results for eight months old children. Quantile functions, raw difference, composition effect, and unexplained difference including support restricted 95% confidence bands. The quantile function lines have been slightly shifted vertically to avoid overlap.

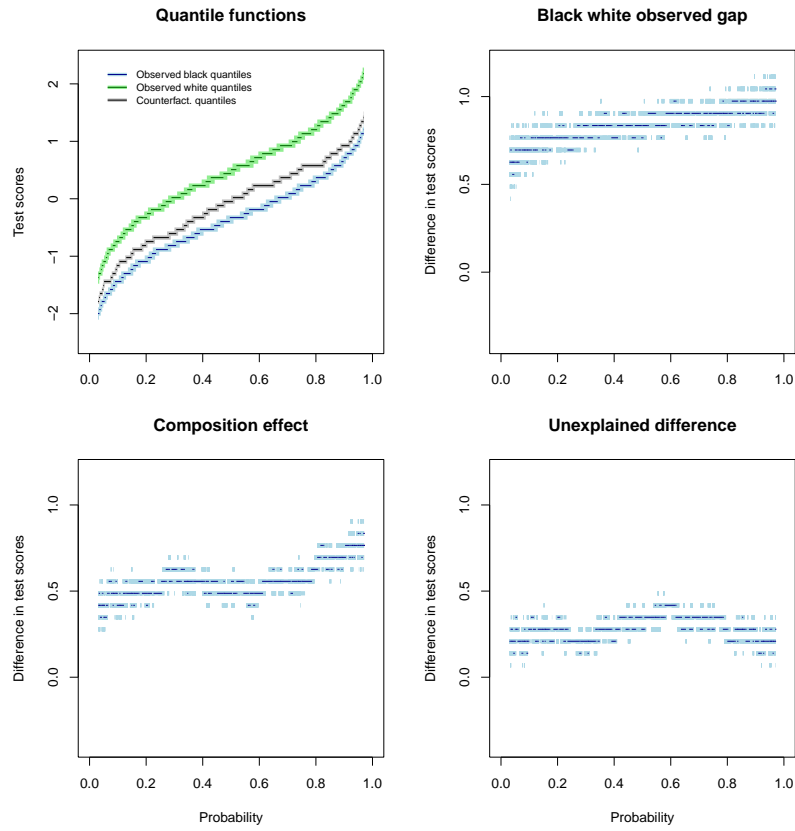


Figure 6: Decomposition of observed racial differences in mental ability of young children; results for seven year old children. Quantile functions, raw difference, composition effect, and unexplained difference including support restricted 95% confidence bands.

tail and is mainly due to unobserved characteristics. While these effects are statistically significant, they are so small in magnitude that they should not worry any policy maker. The composition effect is very small, probably simply because there was no difference to explain to begin with.

The results are completely different for seven years old children. We find a large and statistically significant positive raw black-white gap. A formal test based on the uniform bands rejects the null hypothesis of a zero or a negative racial test score gap at all quantiles. The estimated QE function is increasing in the quantile index ranging from below 0.6 standard deviation units at the lower tail up to over one standard deviation unit at the upper tail of the distribution. The quantile differences at the tails substantially differ from the mean difference of 0.85 standard deviation units reported in Fryer and Levitt (2013). In fact, we can formally reject the null hypothesis of a constant raw test score gap across the distribution because we can not draw a horizontal line at any value of the difference of test scores, which is covered by the confidence band of the QE function at all quantile indexes.

Our decomposition analysis shows that about two third of this gap can be explained by differences in the distribution of observable characteristics. Nevertheless, the remaining unexplained difference is significant, both in economic and in statistical terms. Looking at the QE function, we can see that there is substantial effect heterogeneity along the distribution. Interestingly, the increase in the test score gap at the upper quantiles can be fully explained by differences in background characteristics between black and white children. The resulting unexplained difference is maximized in the center of the distribution. Finally, our simultaneous confidence bands allow for testing several interesting hypothesis' about the whole QE function. For instance, we can reject the null hypothesis that the composition effect and the unexplained difference are zero, negative, or constant at all

quantiles but we cannot reject that they are positive

References

- Angrist, J., Bettinger, E. and Kremer, M. (2006), ‘Long-term educational consequences of secondary school vouchers: Evidence from administrative records in colombia’, *The American Economic Review* **96**(3), 847–862.
- Arcones, M. A. and Giné, E. (1992), ‘On the bootstrap of m-estimators and other statistical functionals’, *Exploring the Limits of Bootstrap*, ed. by R. LePage and L. Billard, Wiley pp. 13–47.
- Baicker, K., Finkelstein, A., Song, J. and Taubman, S. (2014), ‘The impact of medicaid on labor market activity and program participation: Evidence from the oregon health insurance experiment’, *American Economic Review* **104**(5), 322–28.
URL: <http://www.aeaweb.org/articles?id=10.1257/aer.104.5.322>
- Baicker, K., Taubman, S. L., Allen, H. L., Bernstein, M., Gruber, J. H., Newhouse, J. P., Schneider, E. C., Wright, B. J., Zaslavsky, A. M. and Finkelstein, A. N. (2013), ‘The oregon experiment effects of medicaid on clinical outcomes’, *New England Journal of Medicine* **368**(18), 1713–1722. PMID: 23635051.
URL: <http://dx.doi.org/10.1056/NEJMsa1212321>
- Cheng, G., Yu, Z. and Huang, J. Z. (2013), ‘The cluster bootstrap consistency in generalized estimating equations’, *Journal of Multivariate Analysis* **115**, 33–47.
- Chernozhukov, V., Fernandez-Val, I. and Galichon, A. (2009), ‘Improving point and interval estimators of monotone functions by rearrangement’, *Biometrika* p. asp030.
- Chernozhukov, V., Fernández-Val, I. and Melly, B. (2013), ‘Inference on counterfactual distributions’, *Econometrica* **81**(6), 2205–2268.
URL: <http://dx.doi.org/10.3982/ECTA10582>
- Doksum, K. (1974), ‘Empirical probability plots and statistical inference for nonlinear models in the two-sample case’, *The annals of statistics* pp. 267–277.
- Finkelstein, A., Taubman, S., Wright, B., Bernstein, M., Gruber, J., Newhouse, J. P., Allen, H., Baicker, K. and Group, O. H. S. (2012), ‘The oregon health insurance experiment:

- Evidence from the first year*', *The Quarterly Journal of Economics* **127**(3), 1057–1106.
URL: <http://qje.oxfordjournals.org/content/127/3/1057.abstract>
- Foresi, S. and Peracchi, F. (1995), 'The conditional distribution of excess returns: An empirical analysis', *Journal of the American Statistical Association* **90**(430), 451–466.
URL: <http://www.jstor.org/stable/2291056>
- Frydman, H. and Simon, G. (2008), 'Discrete quantile estimation', *Advances and Applications in Statistics* **9**, 177–203.
- Fryer Jr, R. G. and Levitt, S. D. (2013), 'Testing for racial differences in the mental ability of young children', *The American Economic Review* **103**(2), 981–1005.
- Fryer, Roland G., J. and Levitt, S. D. (2013), 'Testing for racial differences in the mental ability of young children', *American Economic Review* **103**(2), 981–1005.
URL: <http://www.aeaweb.org/articles?id=10.1257/aer.103.2.981>
- Galton, F. (1874), 'On a proposed statistical scale', *Nature* **9**, 342–343.
- Imbens, G. W. and Newey, W. K. (2009), 'Identification and estimation of triangular simultaneous equations models without additivity', *Econometrica* **77**(5), 1481–1512.
URL: <http://dx.doi.org/10.3982/ECTA7108>
- Koenker, R. and Bassett Jr, G. (1978), 'Regression quantiles', *Econometrica: journal of the Econometric Society* pp. 33–50.
- Koenker, R. and Xiao, Z. (2002), 'Inference on the quantile regression process', *Econometrica* **70**(4), 1583–1612.
- Kolmogoroff, A. (1933), 'Sulla determinazione empirica di una legge di distribuzione', *Giornale dell'Istituto degli Attuari* **4**, 83–91.
- Kolmogoroff, A. (1941), 'Confidence limits for an unknown distribution function', *The annals of mathematical statistics* **12**(4), 461–463.
- Larocque, D. and Randles, R. H. (2008), 'Confidence intervals for a discrete population median', *The American Statistician* pp. 32–39.
- Ma, Y., Genton, M. G. and Parzen, E. (2011), 'Asymptotic properties of sample quantiles of discrete distributions', *Annals of the Institute of Statistical Mathematics* **63**(2), 227–243.

- Machado, J. A. F. and Silva, J. S. (2005), ‘Quantiles for counts’, *Journal of the American Statistical Association* **100**(472), 1226–1237.
- Praestgaard, J. and Wellner, J. A. (1993), ‘Exchangeably weighted bootstraps of the general empirical process’, *Ann. Probab.* **21**(4), 2053–2086.
URL: <http://dx.doi.org/10.1214/aop/1176989011>
- Scheffe, H. and Tukey, J. W. (1945), ‘Non-parametric estimation. i. validation of order statistics’, *The Annals of Mathematical Statistics* pp. 187–192.
- Sherman, M. and Cessie, S. I. (1997), ‘A comparison between bootstrap methods and generalized estimating equations for correlated outcomes in generalized linear models’, *Communications in Statistics-Simulation and Computation* **26**(3), 901–925.
- Taubman, S. L., Allen, H. L., Wright, B. J., Baicker, K. and Finkelstein, A. N. (2014), ‘Medicaid increases emergency-department use: Evidence from oregon’s health insurance experiment’, *Science* **343**(6168), 263–268.
URL: <http://science.sciencemag.org/content/343/6168/263>
- van der Vaart, A. W. and Wellner, J. A. (1996), *Weak Convergence and Empirical Processes*, Springer Series in Statistics.
- Williams, O. D. and Grizzle, J. E. (1972), ‘Analysis of contingency tables having ordered response categories’, *Journal of the American Statistical Association* **67**(337), 55–63.
- Winkelmann, R. (2006), ‘Reforming health care: Evidence from quantile regressions for counts’, *Journal of Health Economics* **25**(1), 131 – 145.
URL: <http://www.sciencedirect.com/science/article/pii/S0167629605000433>
- Zhang, J. L. and Rubin, D. B. (2003), ‘Estimation of causal effects via principal stratification when some outcomes are truncated by death’, *Journal of Educational and Behavioral Statistics* **28**(4), 353–368.