

**IRISS-C/I**

*An Integrated Research  
Infrastructure in the  
Socio-economic Sciences*

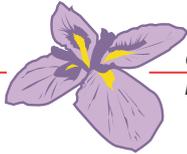
**A Partially Linear Censored Quantile Regression  
Model for Unemployment Duration**

by

Tereza Neocleous

Stephen Portnoy





## **A Partially Linear Censored Quantile Regression Model for Unemployment Duration**

**Tereza Neocleous**

University of Glasgow

**Stephen Portnoy**

University of Illinois

**Abstract** Censored Regression Quantile (CRQ) methods provide a powerful and flexible approach for the analysis of censored survival data when standard linear models are felt to be appropriate. In many cases however, greater flexibility is desired to go beyond the usual multiple regression paradigm. One area of common interest is that of partially linear models, where one (or more) of the explanatory variables are assumed to act on the response through a non-linear function. Here the CRQ approach (Portnoy, 2003) is extended to such partially linear setting. Basic consistency results are presented. A simulation experiment and analysis of unemployment data example justify the use of the partially linear approach over methods based on the Cox proportional hazards regression model and methods not permitting nonlinearity.

**Reference** IRISS Working Paper 2008-07, CEPS/INSTEAD, Differdange, Luxembourg

**URL** <http://ideas.repec.org/p/irs/iriswp/2008-07.html>

The views expressed in this paper are those of the author(s) and do not necessarily reflect views of CEPS/INSTEAD. IRISS Working Papers are not subject to any review process. Errors and omissions are the sole responsibility of the author(s).

# A Partially Linear Censored Quantile Regression Model for Unemployment Duration

Tereza Neocleous<sup>1</sup> and Stephen Portnoy<sup>2</sup>

## Abstract

Censored Regression Quantile (CRQ) methods provide a powerful and flexible approach for the analysis of censored survival data when standard linear models are felt to be appropriate. In many cases however, greater flexibility is desired to go beyond the usual multiple regression paradigm. One area of common interest is that of partially linear models, where one (or more) of the explanatory variables are assumed to act on the response through a non-linear function. Here the CRQ approach (Portnoy (2003)) is extended to such partially linear setting. Basic consistency results are presented. A simulation experiment and analysis of unemployment data example justify the use of the partially linear approach over methods based on the Cox proportional hazards regression model and methods not permitting nonlinearity.

---

<sup>1</sup>Department of Statistics, University of Glasgow  
15 University Gardens , Glasgow G12 8QW, [tereza@stats.gla.ac.uk](mailto:tereza@stats.gla.ac.uk)

<sup>2</sup>Department of Statistics, University of Illinois at Urbana-Champaign  
725 S. Wright St., Champaign IL 61801, [sportnoy@uiuc.edu](mailto:sportnoy@uiuc.edu)

*Keywords:* quantile regression, partially linear models, B-splines  
censored data, unemployment duration

# 1 Introduction

Consider the data which arise from a large scale longitudinal survey to study the durations of spells of unemployment of workers. Exits from unemployment to employment periods are marked and define observed periods of unemployment. Exits from unemployment into states other than employment generate censored values. In this paper we use as an example the German Socio-Economic Panel Survey, where 2214 unemployment durations are observed, of which 55 % are censored. In addition to the unemployment durations, several covariates are observed: gender, marital status, place of residence, age, education and others. The usual approach for the analysis of such data is to express the durations (or log-durations) as a linear model in the covariates and possibly their interactions.

As we discuss below, censored regression quantile methods are especially appropriate when the relationship between outcome and covariates (that is, the parameters or the coefficients of linear regression terms) may be expected to vary with the size of the response, i.e. the conditional quantile, or, more generally because of population heterogeneity. For example, the effect of nationality or gender may be quite different for people with short unemployment durations than for those with longer unemployment spells.

However, even at a fixed quantile, it seems highly unlikely that the effect of age would be strictly linear (even if the data is transformed, say by logarithms). Thus, it is highly desirable to be able to allow the effect of age (and its interaction terms) to be modeled by somewhat nonlinear functions. In this paper, we provide an approach to analysis of such data.

We consider a regression quantile estimator for right censored survival data. Let  $(\mathbf{X}, Y)$  be a random vector with  $\mathbf{X} \in \mathbb{R}^p$  and  $Y$  a real-valued variable.  $\mathbf{X}$  could have discrete or continuous components, with at least one continuous component whose relationship with  $Y$  is nonlinear. For  $\tau \in (0, 1)$  the regression quantile  $Q_{Y|\mathbf{X}}(\tau; \mathbf{x})$  of  $Y$  given  $\mathbf{X} = \mathbf{x}$  satisfies

$$P(Y \leq Q_{Y|\mathbf{X}}(\tau; \mathbf{x}) | \mathbf{X} = \mathbf{x}) = \tau.$$

Assuming that  $n$  independent pairs  $(Y_i, \mathbf{X}_i)$  are observed, and that the relationship between

$Y$  and  $\mathbf{X}$  is linear, i.e.

$$Q_{Y_i}(\tau|\mathbf{x}_i) = \mathbf{x}_i^\top \boldsymbol{\beta}(\tau), \quad (1)$$

the  $\tau$ th regression quantile coefficient,  $\hat{\boldsymbol{\beta}}(\tau)$ , and hence the regression quantile  $\hat{Q}_{Y|\mathbf{X}}(\tau; \mathbf{x})$ , can be obtained as the solution of

$$\min_{\mathbf{b} \in \mathbb{R}^p} \sum_{i=1}^n \rho_\tau(y_i - \mathbf{x}_i^\top \mathbf{b})$$

where  $\rho_\tau(u) = u(\tau - I(u < 0))$  (see Koenker (2005) for details). With survival times it is often the case that  $Y$  is not observed, and that instead one observes only the minimum of  $Y$  and a censoring variable  $C$ . Suppose that  $n$  independent triples  $\{(\mathbf{X}_i, Z_i, \Delta_i), i = 1, \dots, n\}$  are observed, with  $Z_i = \min(Y_i, C_i)$  and  $\Delta_i = I(Y_i \leq C_i)$ . We are interested in estimating  $Q_{Y|\mathbf{X}}(\tau; \mathbf{x})$  when  $Y$  and  $C$  are conditionally independent given  $\mathbf{X}$ , and when  $Y$  varies linearly with most components of  $\mathbf{X}$  but nonlinearly with at least one component of  $\mathbf{X}$ .

Under the linear models paradigm a quantile regression approach is especially useful in survival analysis, as it interprets the covariate effect on survival times with flexibility not always achievable under the global assumptions like those of the Cox model. Koenker and Geling (2001) introduced a quantile regression approach to survival analysis by means of a transformation of the survival times. For instance, when the log-transformation is used, quantile regression corresponds to the accelerated failure time model, in which  $\log Y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + u_i$  and the hazard rate is given by

$$h_i(y|\mathbf{x}) = h_0(y \exp(-\mathbf{x}_i^\top \boldsymbol{\beta})) \exp(-\mathbf{x}_i^\top \boldsymbol{\beta}).$$

Moreover, if the  $u_i$  are i.i.d. with extreme value distribution  $F(u) = 1 - \exp(-\exp(u))$ , this corresponds to the Cox proportional hazards model with Weibull baseline hazard, and the linear quantile regression model for the log-survival times agrees with the Cox model for accelerated failure time. Otherwise the Cox model specifies a parametric model for the survival distribution, while quantile regression permits rather general heterogeneity (subject to the use of linear models). The proportional hazards model is the most popular method

for analyzing right-censored survival data, but in recent years there have been advances in quantile regression methods that offer an alternative to the Cox approach.

The earliest proposed estimator for censored quantile regression assumed fixed censoring (Powell (1986)). Subsequent research either assumed fixed censoring or independence between  $Y$  and  $C$ , *e.g.* Buchinsky and Buchinsky and Hahn (1998), Honore et al (2002), and Chernozhukov and Hong (2002).

The independence assumption was relaxed in Portnoy (2003), where conditional independence of  $Y$  and  $C$  given  $\boldsymbol{x}$  is assumed, and a “reweighting-to-the-right” (Efron (1967)) scheme is employed to compute the conditional quantiles. The Portnoy (2003) method is of particular interest, as it essentially extends the Kaplan-Meier estimator to the regression setting. A similar generalization of the Nelson-Aalen estimator was also recently proposed by Peng and Huang (2008). The models developed in the rest of this paper are based on the Portnoy estimator.

The Portnoy CRQ model assumes conditional independence between  $Y_i$  and  $C_i$  given  $\boldsymbol{x}_i$ . The approach is based on a recursive pivoting algorithm for random censoring, whose solution reduces to the Kaplan-Meier estimator in the one-sample case. The algorithm iteratively computes the entire conditional quantile function for  $\tau \in (0, 1)$ , stopping at a value of  $\tau$  for which all observations remaining above the current conditional quantile function are censored. Note that this differs from the usual quantile regression methods that compute the conditional quantile at a fixed  $\tau$ . If, for instance, the median is required, the pivoting algorithm of Portnoy (2003) will compute all quantiles up to the 50th in order to obtain the median.

In what follows, we present a modification of the pivoting algorithm with a generalization permitting nonlinear response to one (or more) covariates (as a “partially linear” model). Section 2 presents a grid algorithm as a computationally effective method for fitting such models based on generally available regression quantile programs. Section 3 examines the asymptotic properties of the partially linear CRQ estimator. Simulation experiments are

statistically analyzed in Section 4 to evaluate the performance of the approach. A study of unemployment duration data is presented in Section 5 to show the value of the use of the partially linear censored regression model.

## 2 Grid algorithm for linear CRQ

A slightly modified version of the Portnoy (2003) CRQ pivoting algorithm, evaluating the linear regression quantiles of (1) on a grid of  $\tau$  values is presented here. This algorithm iteratively computes the conditional quantiles from lowest to highest. Suppose that at the starting value  $t_1$  of  $\tau \in (0, 1)$  there are no censored observations below the  $t_1$ th quantile, so that the quantile coefficient  $\hat{\beta}(t_1)$  is estimated using the usual quantile regression algorithm minimizing  $\sum_{i=1}^n \rho_{t_1}(y_i - \mathbf{x}_i^T \mathbf{b})$  with respect to  $\mathbf{b}$ . The corresponding quantile hyperplane  $\mathbf{x}_i^T \hat{\beta}(t_1)$  will then have proportion  $t_1$  of the data below it and  $(1 - t_1)$  above. We say that observations for which  $Y_i \leq \mathbf{x}_i^T \hat{\beta}(t_1)$  are *crossed* by the  $t_1$ th quantile. As the value of  $\tau$  increases, censored observations may also get crossed. When the  $i$ th censored observation is crossed, the algorithm splits it to two parts according to a weighting scheme: a part that is observed at  $C_i$  and a part at infinity. If the  $i$ th censored point  $C_i$  is crossed for the first time at  $\tau = \tau_i$ , it will receive weight  $\hat{w}_i(\tau) = (\tau - \tau_i)/(1 - \tau_i)$  for all  $\tau > \tau_i$ . This weight is updated every time  $\tau$  increases. With weights for all crossed censored observations computed, weighted quantile regression is performed to obtain the regression coefficients at the current value of  $\tau$ . More details on the weights of crossed observations and on the weighted quantile regression performed are given below.

### Algorithm

1. Choose gridpoints  $t_1, \dots, t_M$  covering the set  $\varepsilon \leq \tau \leq 1 - \varepsilon$ . Starting with the gridpoint  $t_1$  compute the initial quantile function  $\hat{\beta}(t_1)$  for  $1 \leq \tau \leq t_1$  using the uncensored quantile regression algorithm. This assumes that the initial regression quantile,  $\hat{\beta}(t_1)$  determines a hyperplane that lies below all censored points, which is reasonable, since censored observations below all data are non-informative and can be deleted without

changing the estimation.

2. Suppose that the quantiles  $\hat{\boldsymbol{\beta}}(t_l), 1 \leq l \leq k$  have been computed by minimizing over  $\mathbf{b} \in \mathbb{R}^p$  the objective function

$$\sum_{i=1}^n \left\{ \begin{aligned} &\Delta_i \rho_{t_{k+1}}(Y_i - \mathbf{x}_i^\top \mathbf{b}) \\ &+ (1 - \Delta_i) [\hat{w}_i(t_{k+1}, \boldsymbol{\beta}) \rho_{t_{k+1}}(C_i - \mathbf{x}_i^\top \mathbf{b}) \\ &+ (1 - \hat{w}_i(t_{k+1}, \boldsymbol{\beta})) \rho_{t_{k+1}}(Y^* - \mathbf{x}_i^\top \mathbf{b})] \end{aligned} \right\}$$

where  $Y^*$  is a sufficiently large value so that  $Y^* > \mathbf{x}_i^\top \mathbf{b}$  for all  $\mathbf{x}_i^\top \mathbf{b}$  from the data.  $Y^*$  will be referred to as “point at infinity”.

3. In the step from  $t_k$  to  $t_{k+1}$  some censored observations that were not previously crossed might get crossed. For those observations  $C_i > \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}(t_k)$  and  $C_i \leq \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}(t_{k+1})$ . They are then given weights  $\hat{w}_i(\tau) = (\tau - \tau_i)/(1 - \tau_i)$  with  $\tau_i(\hat{\boldsymbol{\beta}}) = t_k$  with the rest of the weight going to the point at infinity,  $Y^*$ . In addition, updated weights are computed for the already crossed observations according to the same formula. With all the weights defined, a usual weighted quantile regression is performed.
4. The algorithm stops either at the last grid point,  $t_M$ , or at some point  $t_e$  when only non-reweighted censored observations remain above the current solution,  $\mathbf{x}_i^\top \hat{\boldsymbol{\beta}}(t_e)$ .

The main advantage of using the grid modification of the pivoting algorithm is computational. For large sample sizes the pivoting algorithm computes solutions at a high number of  $\tau$ -values. With the grid algorithm the number of  $\tau$ -values at which the solution is obtained can be reduced, with substantial savings in computational time required for the iterative process. The grid algorithm is outlined above for a linear CRQ model, for which asymptotic results are given in Vanden Branden (2005) and Neocleous et al (2006). In what follows the algorithm is applied within the framework of partially linear models.

### 3 The partially linear estimator and its large sample properties

The partially linear CRQ model combines semiparametric estimation for censored data with quantile regression techniques, and uses B-splines for the estimation of the nonlinear term. Consider first the uncensored fully nonlinear model  $y_i = g_\tau(x_i) + e_i$ , where the  $e_i$  are independent random errors with  $\tau$ th quantile equal to zero. Following the notation in Schumaker (1981), let

$$\pi(s) = (B_1(s), B_2(s), \dots, B_{k_n'+d+1}(s))^\top$$

be the set of B-spline basis functions with given knots  $\Delta = \{x_i\}_0^{k_n'}$  with number of spline knots  $k_n'$  and order of splines  $d + 1$ . Then the estimated  $\tau$ th quantile function  $\hat{g}_{n\tau}(s) = \pi(s)^\top \hat{\boldsymbol{\theta}}_n$ , where  $\hat{\boldsymbol{\theta}} \in \mathbb{R}^{k_n'+d+1}$ , is a solution of

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^{k_n'+d+1}} \sum_i \rho(y_i - \pi(x_i)^\top \boldsymbol{\theta}).$$

Once the spline knots are selected and the spline bases computed, the problem is reduced to a linear quantile regression with  $(k_n' + d + 1)$  parameters. It was shown, *e.g.* in He and Shi (1994, 1996) that if  $g_\tau$  is smooth with bounded  $r$ th derivative, and  $k_n'$  is of order  $n^{1/(2r+1)}$ , under some mild conditions the spline estimate  $\hat{g}_{n\tau}(s)$  converges to  $g_\tau(s)$  at the optimal nonparametric rate of  $\mathcal{O}_p(n^{-2r/(2r+1)})$ . In what follows we discuss the use of a B-spline estimator in a censored regression quantile setting.

Assuming the data  $\mathbf{x}_i = (\mathbf{x}_{1i}, \mathbf{x}_{2i})$ ,  $i = 1, \dots, n$ , come from a model with

$$Q_{Y_i}(\tau | \mathbf{x}_i) = \mathbf{x}_{1i}^\top \boldsymbol{\theta}_1(\tau) + g_\tau(\mathbf{x}_{2i}), \quad (2)$$

the estimated quantiles will be of the form

$$\hat{Q}_{Y_i}(\tau | \mathbf{x}_i) = \mathbf{x}_{1i}^\top \hat{\boldsymbol{\theta}}_1(\tau) + \pi(\mathbf{x}_{2i})^\top \hat{\boldsymbol{\theta}}_2(\tau), \quad (3)$$

where  $g_\tau$  is approximated by a linear combination of B-splines.

Let  $\boldsymbol{\beta} = (\boldsymbol{\theta}_1^\top, \boldsymbol{\theta}_2^\top)^\top$ . Without loss of generality, we assume that the support of  $g(s)$  is  $s \in [0, 1]$ . Let  $\pi(s) = (\pi_1(s), \pi_2(s), \dots, \pi_{k_n'+d+1}(s))^\top$  be the B-spline basis of order  $d$  with

$k'_n$  knots. Let  $k_n = k'_n + d + 1$  and define  $R_i(\tau) = \pi(\mathbf{x}_{2i})^\top \theta_2(\tau) - g_\tau(\mathbf{x}_{2i})$ . Then at the  $k$ th step of the CRQ grid algorithm

the estimated  $t_{k+1}$ th quantile is  $\mathbf{x}_{1i}^\top \hat{\theta}_1(t_{k+1}) + \pi(\mathbf{x}_{2i})^\top \hat{\theta}_2(t_{k+1})$ . This linearity in  $\beta$  allows current theoretical approaches to be generalized to the case of  $\beta$  of increasing dimension (at the same rate as  $k_n$ ). For a grid of  $M$   $\tau$ -values the CRQ estimator is  $\hat{\beta} = (\hat{\beta}(t_1)^\top, \hat{\beta}(t_2)^\top, \dots, \hat{\beta}(t_M)^\top)^\top \in \mathbb{R}^{Mp}$  and the following result holds.

**Theorem 3.1** *Let  $\hat{\beta} \in \mathbb{R}^{Mp}$ , be the censored regression quantile estimator for the model specified in (1) on a grid  $\varepsilon \leq t_1 < t_2 < \dots < t_M \leq 1 - \varepsilon$ . Let  $\beta^*$  be the true unknown censored regression quantile along the same grid,  $t_{k+1} - t_k \equiv g_n = n^{-\kappa}$  and  $p = \mathcal{O}(n^\gamma)$  where  $\gamma$  and  $\kappa$  satisfy one of (4), (5) and (6):*

$$0 < \kappa < 1/6, \quad 0 < \gamma < \kappa \tag{4}$$

$$1/6 < \kappa < 1/4, \quad 0 < \gamma < 1/4 \tag{5}$$

$$1/4 < \kappa < 1/3, \quad 0 < \gamma < (1 - 3\kappa)/2. \tag{6}$$

*Under Assumptions (I), (F), (X) and (XX) given in the Appendix,*

$$\|\hat{\beta} - \beta^*\|^2 = \mathcal{O}_p(n^{\kappa+\gamma-1}).$$

For the partially linear CRQ model with B-spline estimation of the nonlinear part, the following corollary holds.

**Corollary 3.2** *Let  $\hat{\beta} = (\hat{\theta}_1^\top, \hat{\theta}_2^\top)^\top \in \mathbb{R}^{Mp}$  be the censored regression quantile grid estimator of  $\beta^* = (\theta_1^{*\top}, \theta_2^{*\top})$ , where  $\pi(x_2)^\top \theta_2^*$  estimates  $g(x_2)$  in the model specified in (2). Under the assumptions and notation of Theorem 3.1, with the added condition*

*(G)  $g_\tau(s)$  has bounded  $r$ th derivative for  $r \geq 3$  for all  $\tau$ ,*

$$\|\hat{\theta}_1 - \theta_1^*\|^2 = \mathcal{O}_p(n^{\kappa+\gamma-1}).$$

Corollary 3.2 can be proved by combining B-spline approximation rates and Theorem 3.1. This result is most useful in applications where the effect of interest, *e.g.* treatment effect, is to be estimated in the presence of some additional nonlinear covariate.

## 4 Simulation study

To examine the finite sample performance of the partially linear CRQ estimator, we conducted a simulation experiment in which the censored response is linear in one covariate and non-linear in another covariate. Event times were generated for  $i = 1, \dots, n$  from the model

$$Y_i = \beta_0 + \beta_1 x_{1i} + \frac{10e_{1i}}{1 + \exp(6 - 0.5x_{2i})}$$

and censoring times from the model (Configuration 1)

$$C_i = \beta_0 + \beta_1 x_{1i} + \frac{10e_{2i}}{1 + \exp(5 - 0.5x_{2i})}$$

for roughly 20% censoring, or (Configuration 2)

$$C_i = \beta_0 + \beta_1 x_{1i} + \frac{10e_{2i}}{1 + \exp(4 - x_{2i})} - 0.2x_1^2$$

for roughly 40% censoring. Parameter values were  $\beta_0 = 1$  and  $\beta_1 = 3$ , and the  $x_{1i}$  were generated as iid  $U(0, 5)$ , the  $x_{2i}$  as iid  $U(0, 25)$ , and  $e_{1i}$  and  $e_{2i}$  as iid  $N(1, 0.01)$ . The scatterplot in Figure 1 shows the censoring mechanism for Configuration 1 and sample size  $n = 500$ . Four different models were fitted to the data: one with linear term in  $x_2$  and three with spline terms of order 2, 3 and 4 (piecewise linear, quadratic and cubic) in  $x_2$ . Knots at the quartiles of  $x_2$  were used in the spline models for Configuration 1, while for Configuration 2 two additional sets of knots were considered. In each case bootstrap confidence intervals were computed with  $b = 500$  bootstrap replications.

Tables 1 and 2 report average bias, median absolute error, root mean square error, empirical coverage probability (95% nominal coverage) and mean confidence interval length for the slope of  $x_1$  evaluated at  $\tau = 0.50$  and  $0.75$  (similar results were obtained for  $\tau = 0.25$ ) for Configuration 1. In all cases the partially linear model outperforms its linear equivalent.

The difference between the three spline orders used is less clear, with some evidence that the quadratic spline works best. This is also supported by Figure 2, in which the quadratic spline term appears to give the best fit for the nonlinear term.

The effect of knot selection and placement is further investigated in the simulation study of Configuration 2, in which fitted spline models have knots at (a) the 33rd and 66th quantile of  $x_2$ , (b) the quartiles of  $x_2$ , and (c) the 20th, 40th, 60th and 80th quantiles of  $x_2$ . Tables 3 and 4 show the performance of various models fitted for Configuration 2. It can be seen that again the spline models perform better than the linear model, while three knots are in general better than just two. The difference between three and four knots is less clear, as it appears that three knots are better for quadratic spline models, and four knots are better for piecewise linear and cubic spline models.

Finally, Table 5 reports bias, root mean square error and median absolute error for the estimation of the nonlinear term in Configuration 2. The quadratic spline with three knots appears to be performing better than other spline models in terms of root mean square error. Differences in bias are less obvious.

## 5 Application to unemployment duration

We illustrate the usefulness of the partially linear CRQ model with an application to administrative unemployment data from the German Socio-Economic Panel Survey, a longitudinal survey of private households in Germany covering topics such as income, employment, education and health. We focus on a subset of the data covering the period 1992-2004. The response variable of interest,  $Y$ , is the duration in months of the latest unemployment spell in the respondent's work history.

We restrict our attention to males with German nationality (as both nationality and gender were found to be significant in preliminary analyses) and we explore the effect of age and marital status on unemployment duration. Exits from unemployment to full- or part-time employment were considered observed while all other exits were considered as censored observations. Excluding observations with missing data, this gave a sample size

of 2214 records with 55% censoring. Of these 2214 individuals, 42% were married. The median age for married respondents was 47.42 and for single 26.17.

The CRQ model

$$Q_{\log(Y_i)}(\tau | \mathbf{x}) = \beta_0(\tau) + \beta_1(\tau) \times \text{married} + \boldsymbol{\theta}(\tau)^\top \pi(\text{age}) \quad (7)$$

was considered and quantiles up to the 60th were estimated. In particular, a quadratic spline term with knots at the quartiles of age was fitted. This provides a smooth 5-parameter fit to the age effect. All but one of the five coefficients were significant (at some  $\tau$ -values), and so it is clear that the age effect requires more than a linear term.

Plots of  $\hat{\beta}(\tau)$ , the estimated quantile coefficients for the intercept and marital status, against  $\tau$  are shown in Figure 3. The coefficients tend to be smaller in absolute value for short term unemployment and larger for long term unemployment.

Marriage has a strong negative effect on unemployment duration, independent of age (the relevant interaction terms were not significant). The estimated median coefficient representing the difference in log-duration between a single and a married German male is -0.8244 (confidence interval of (-1.1649,-0.4838)), i.e. median unemployment duration for a married respondent is 0.4385 times that of a single respondent of the same age. The size of the marriage effect is similar in all but the lowest quantiles of unemployment duration.

Plots of the estimated median unemployment duration against age are shown in Figure 4 separately for single and married German males. Pointwise bootstrap confidence intervals are also shown. The age ranges plotted reflect the different age distributions for married and single groups. For married males over 50, censoring exceeds 80%, thus we restrict attention for the married group to the “reliable estimation” age range (31.42,50.00) corresponding to the 10th age percentile and the age with 81% censoring above it. For single males the age range plotted is (19.67,47.17) corresponding to the 10th and 90th age percentiles. In the singles age distribution, 80% of the observations over age 47.17 are censored.

From Figure 4, it is clear that the age effect on unemployment duration is quite nonlinear (at least for single men), with age being beneficial at very low ages ( $< 25$ ) and rather

detrimental (for both single and married men) at higher ages (as might be expected). The quantile analysis in Figure 5 presents perhaps a more surprising result. For quantiles below  $\tau = .3$  (shorter unemployment durations), the effect is rather independent of age. This is not unexpected, as those who are readily re-employable do well at any age. However, for higher quantiles, the detrimental effect of age seems to increase rapidly for men in the range 30 - 50 years. The rather substantial increase in difficulty to obtain employment for older men who are not so readily re-employable would seem to call for some explanation (economic, psychological, or sociological).

Plots such as those in Figures 3-5 are useful in identifying departures from linearity. We advocate exploring the nonlinearity of each continuous covariate before attempting to fit linear coefficients as a way to detect patterns and improve the overall fit of the model. In addition, fitting a CRQ model can highlight differences in the covariate effects for long and short-term durations, something that is not picked up by the proportional hazards model.

## 6 Concluding remarks

In the preceding sections we proposed the use of a partially linear model for censored regression quantiles as a useful extension to the standard linear regression techniques for survival data. The partially linear model was shown to be consistent and its use was illustrated by a data example and simulations. Quartile knots were used for the B-spline estimation of nonlinear terms and the quadratic spline gave satisfactory quantile estimates in the empirical example and simulations. Higher order spline terms did not show much improvement in estimation.

The censored regression quantile estimator is robust and flexible enough to highlight aspects of the data that the most common survival analysis techniques might overlook. Incorporating a nonlinear part adds even more flexibility to the model, allowing for more accurate estimation of parameters of interest, like quantile treatment effects. Censored regression quantiles and the semiparametric model proposed here are tools for capturing subtle aspects of the data and can be used in conjunction with other techniques for more

comprehensive exploration of censored data.

As in every semiparametric model, the use of B-splines raises the question of knot selection. In this work the spline knots were chosen at fixed quantiles of the nonlinear variable. As long as the knot selection is not data-driven (*e.g.* equally spaced knots or quantile knots, perhaps depending on the sample size  $n$ ), the asymptotic theory of B-splines applies directly (and consistency follows by Theorem 3.1 if the number of knots increases with  $n$  appropriately). Asymptotic results are not currently available if knot selection is data-driven. In practice fixing knots at specified quantiles of the  $x$ -variable is a simple and convenient solution for small to medium-sized datasets, and it is not likely that data-driven methods can offer much improvement here. However, in general it is also desirable to have a method for optimal knot selection and placement depending on the data. Such methods have been proposed by a number of authors. For instance, Koenker et al (1994) use a roughness penalty for quantile smoothing splines, and Doksum and Koo (2000) propose a method for stepwise knot addition and deletion using modified AIC and BIC for nonparametric quantile regression with regression splines. Further work along such lines would be useful for larger data sets.

## Acknowledgments

This research was partially supported by National Science Foundation Grant DMS06-04229 and by the European Commission under the 6th Framework Programme's Research Infrastructures Action (Trans-national Access contract RITA 026040) hosted by IRISS-C/I at CEPS/INSTEAD, Differdange (Luxembourg). The data used in this publication were made available to us by the German Socio-Economic Panel Study (SOEP) at the German Institute for Economic Research (DIW), Berlin.

## References

Buchinsky M, Hahn JY (1998) An alternative estimator for the censored quantile regression model. *Econometrica* 66:653–671

- Chernozhukov V, Hong H (2002) Three-step censored quantile regression and extramarital affairs. *J Amer Statist Assoc* 97:872–882
- Doksum K, Koo J (2000) On spline estimators and prediction intervals in nonparametric regression. *Computational Statistics and Data Analysis* 35:67–82
- Efron B (1967) The two-sample problem with censored data. vol 4, pp 831–853, proceedings of the Fifth Berkeley Symposium
- He X, Shao Q (2000) On parameters of increasing dimensions. *J Multivariate Anal* 73:120–135
- He X, Shi P (1994) Convergence rate of B-spline estimators of nonparametric conditional quantile functions. *Nonparametric Statistics* 3:299–308
- He X, Shi P (1996) Bivariate tensor-product B-splines in a partly linear model. *J Multivariate Anal* 58:162–181
- Honore B, Khan S, Powell JL (2002) Quantile regression under random censoring. *J Econometrics* 109:67–105
- Koenker R (2005) *Quantile Regression*. Cambridge University Press
- Koenker R, Geling O (2001) Reappraising medfly longevity: a quantile regression survival analysis. *J Amer Statist Assoc* 96:458–468
- Koenker R, Ng P, Portnoy S (1994) Quantile smoothing splines. *Biometrika* 81:673–680
- Neocleous T, Vanden Branden K, Portnoy S (2006) Correction to: “Censored Regression Quantiles”. *J Amer Statist Assoc* 101:860–861
- Peng L, Huang Y (2008) Survival analysis with quantile regression models. *Journal of the American Statistical Association* (103):637–649
- Portnoy S (2003) Censored regression quantiles. *J Amer Statist Assoc* 98:1001–1012

Powell JL (1986) Censored regression quantiles. *J Econometrics* 32:143–155

Vanden Branden K (2005) Robust methods for high-dimensional data, and a theoretical study of depth-related estimators. PhD thesis, Katholieke Universiteit Leuven

## Appendix: Proof of Theorem 3.1

The conditions for the main result (Theorem 3.1) are as follows:

(I)  $Y$  and  $C$  are conditionally independent given  $\mathbf{x}$

(F) For  $0 < \varepsilon < 1$ , there exist constants  $a_j, b_j, c_j$  with  $a_j > 0$  and  $b_j < \infty$  for  $j = 1, 2, 3$

such that

$$\begin{aligned} a_1 \leq f_{Y_i}(y) \leq b_1 & & |f'_{Y_i}(y)| \leq c_1 \\ a_2 \leq \tilde{f}_{Y_i}(u) \leq b_2 & & |\tilde{f}'_{Y_i}(u)| \leq c_2 \\ a_3 \leq \tilde{f}_{C_i}(v) \leq b_3 & & |\tilde{f}'_{C_i}(v)| \leq c_3 \end{aligned}$$

uniformly for  $\varepsilon \leq F_{Y_i}(y) \leq 1 - \varepsilon$ ,  $\varepsilon \leq \tilde{F}_{Y_i}(u) \leq 1 - \varepsilon$  and  $\varepsilon \leq \tilde{F}_{C_i}(v) \leq 1 - \varepsilon$  and uniformly in  $i = 1, \dots, n$ .

(X)  $\max_{1 \leq i \leq n} \|\mathbf{x}_i\| = \mathcal{O}(p)$ .

(XX) The matrix  $\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$  is positive definite.

Theorem 3.1 makes use of the theory of He and Shao (2000) on the asymptotics of M-estimators when the parameter dimension increases with  $n$ . Briefly, this is outlined as follows. Let  $\hat{\boldsymbol{\beta}}_n \in \mathbb{R}^m$  be the M-estimator for minimizing  $\sum_{i=1}^n \zeta(\mathbf{z}_i, \boldsymbol{\beta})$  for some data set  $\{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n\}$  with  $\mathbf{z}_i \in \mathbb{R}^{p+1}$  for  $i = 1, 2, \dots, n$ ; and for some objective kernel  $\zeta(\mathbf{z}_i, \boldsymbol{\beta})$ . If the objective function is convex in  $\boldsymbol{\beta}$ , and if  $\zeta(\mathbf{z}, \boldsymbol{\beta})$  is differentiable with respect to  $\boldsymbol{\beta}$ , except at finitely many points, with derivative  $\Psi(\mathbf{z}, \boldsymbol{\beta})$ , then Theorem 2.1 of He and Shao (2000) states that under certain conditions,  $\|\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}^*\|^2 = \mathcal{O}_p(m/n)$  where  $\boldsymbol{\beta}^*$  is the solution to  $\sum_{i=1}^n E_{\boldsymbol{\beta}} \Psi(\mathbf{z}_i, \boldsymbol{\beta}) = 0$ . For the CRQ grid estimator the increasing dimension is  $m = Mp$ , where  $M$  is the number of grid points. Let  $p = \mathcal{O}(n^\gamma)$  for some  $\gamma > 0$ . Equivalently,  $p \leq cn^\gamma$  for some constant  $c$ . Define  $\Psi_k(\mathbf{x}_i, \boldsymbol{\beta}) = \mathbf{x}_i \{ \Delta_i(I(Y_i < \mathbf{x}_i^\top \boldsymbol{\beta}(t_k))) + (1 - \Delta_i)(w_i(\boldsymbol{\beta}, t_k)I(C_i <$

$\mathbf{x}_i^\top \boldsymbol{\beta}(t_k) - t_k\}$ ,

$$\eta_i(\theta, \boldsymbol{\beta}) = \Psi(\mathbf{x}_i, \boldsymbol{\beta}) - \Psi(\mathbf{x}_i, \boldsymbol{\beta}) - E(\Psi(\mathbf{x}_i, \theta) - \Psi(\mathbf{x}_i, \boldsymbol{\beta}))$$

and  $S_m = \{\alpha \in \mathbb{R}^m : \|\alpha\| = 1\}$ . Then

$$\Psi(\mathbf{x}_i, \boldsymbol{\beta}) = (\Psi_1(\mathbf{x}_i, \boldsymbol{\beta})^\top, \Psi_2(\mathbf{x}_i, \boldsymbol{\beta})^\top, \dots, \Psi_M(\mathbf{x}_i, \boldsymbol{\beta})^\top)^\top \in \mathbb{R}^m.$$

The result also relies on the following two lemmas, which have been shown in the case of fixed  $p$  by Vanden Branden (2005). Here the result is extended to the case of  $p$  growing with  $n$ . Lemma 1 permits restricting the proof to monotone functions  $\mathbf{x}^\top \boldsymbol{\beta}(\tau)$  on the grid. Lemma 2 shows that  $\tau_i(\boldsymbol{\beta})$  and  $\tau_i(\boldsymbol{\beta}^*)$  are close on the set of slopes  $\boldsymbol{\beta}$  considered.

**Lemma 1** *For every  $B > 0$ ,  $\exists n_0$  such that for  $n \geq n_0$  the set*

$$\left\{ \boldsymbol{\beta} \in \mathbb{R}^m : \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\| \leq B \left( \frac{m}{n} \right)^{1/2} \right\}$$

*is contained in the set of all monotonic functions on the grid  $\varepsilon \leq t_1 < t_2 < \dots \leq t_M \leq 1 - \varepsilon$  for some  $\varepsilon > 0$ . Here  $t_k - t_{k-1} = g_n = n^{-\kappa}$ ,  $p \leq cn^\gamma$  for some  $c > 0$ , and  $m \leq p/g_n$ , with  $\gamma \leq \frac{1}{2} - \frac{3\kappa}{2}$ ,  $\kappa > 0$ .*

**Lemma 2** *Let  $\tau_i(\boldsymbol{\beta})$  be the gridpoint at which  $\boldsymbol{\beta}$  crosses  $C_i$ , and let  $\tau_i(\boldsymbol{\beta}^*)$  be the unknown gridpoint at which the true regression quantile  $\boldsymbol{\beta}^*$  crosses the same observation. It then holds that*

$$|\tau_i(\boldsymbol{\beta}) - \tau_i(\boldsymbol{\beta}^*)| = \mathcal{O}(T(n, m))$$

*on the set  $\{\boldsymbol{\beta} : \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\| \leq B(m/n)^{1/2}\}$  with*

$$T(n, m) = \max(Bc^{1/2}p^{1/2}(m/n)^{1/2}, 2g_n) = \max(Bcn^{\kappa+\gamma-1/2}, 2n^{-\kappa}).$$

Proofs of Lemmas 1 and 2 are straightforward generalizations of those in Vanden Branden (2005).

**Proof of Theorem 3.1.** It is sufficient to verify the following conditions of He and Shao (2000).

$$(C0) \quad \|\sum_{i=1}^n \Psi(\mathbf{x}_i, \hat{\boldsymbol{\beta}}_n)\| = o_p(n^{1/2}).$$

(C1) There exists a  $C$  and  $r \in (0, 2]$  such that

$$\max_{i \leq n} E_{\boldsymbol{\beta}} \sup_{\boldsymbol{\theta}: \|\boldsymbol{\theta} - \boldsymbol{\beta}\| \leq d} \|\eta_i(\boldsymbol{\theta}, \boldsymbol{\beta})\|^2 \leq n^C d^r$$

for  $0 < d \leq 1$ .

$$(C2) \quad \|\sum_{i=1}^n \Psi(\mathbf{x}_i, \boldsymbol{\beta}^*)\| = \mathcal{O}_p(nm)^{1/2} \text{ or } \sum_{i=1}^n E \|\Psi(\mathbf{x}_i, \boldsymbol{\beta}^*)\|^2 = \mathcal{O}(nm).$$

(C3) There exists a sequence of  $(m \times m)$  matrices  $D_n$  with  $\liminf_{n \rightarrow \infty} \lambda_{\min}(D_n) > 0$  (where  $\lambda_{\min}$  denotes the minimum eigenvalue) such that for any  $B > 0$  and uniformly in  $\boldsymbol{\alpha} \in S_m$

$$\sup_{\|\boldsymbol{\beta} - \boldsymbol{\beta}^*\| \leq B(\frac{m}{n})^{1/2}} |\boldsymbol{\alpha}^\top \sum_{i=1}^n E_{\boldsymbol{\beta}^*} (\Psi(\mathbf{x}_i, \boldsymbol{\beta}) - \Psi(\mathbf{x}_i, \boldsymbol{\beta}^*)) - n \boldsymbol{\alpha}^\top D_n (\boldsymbol{\beta} - \boldsymbol{\beta}^*)| = o(n^{1/2}).$$

(C4) There exists a sequence  $A(n, m) = o(n/\log n)$  for which

$$\sup_{\boldsymbol{\beta}: \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\| \leq B(\frac{m}{n})^{1/2}} \sum_{i=1}^n E_{\boldsymbol{\beta}} |\boldsymbol{\alpha}^\top \eta_i(\boldsymbol{\beta}, \boldsymbol{\beta}^*)|^2 = \mathcal{O}(A(n, m))$$

for any  $\boldsymbol{\alpha} \in S_m$ , and  $B > 0$ .

$$(C5) \quad \sup_{\boldsymbol{\alpha} \in S_m} \sup_{\boldsymbol{\beta}: \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\| \leq B(\frac{m}{n})^{1/2}} \sum_{i=1}^n (\boldsymbol{\alpha}^\top \eta_i(\boldsymbol{\beta}, \boldsymbol{\beta}^*))^2 = \mathcal{O}_p(A(n, m)) \text{ for any } B > 0.$$

(C0) follows from the gradient conditions by noting that

$$\|\Psi(\hat{\boldsymbol{\beta}})\|^2 = \mathcal{O}_P(M \max_{1 \leq k \leq M} \|\Psi_k(\hat{\boldsymbol{\beta}}(t_k))\|^2)$$

and

$$\|\Psi_k(\hat{\boldsymbol{\beta}})\| = \mathcal{O}_P(\sqrt{p \log n} \max \|\mathbf{x}_i\|).$$

Thus

$$\|\Psi(\hat{\boldsymbol{\beta}})\| = \mathcal{O}_P(p \sqrt{M \log n}) = \mathcal{O}_P(n^{\kappa/2 + \gamma} (\log n)^{1/2}).$$

This is  $o_p(n^{1/2})$ , provided that  $\kappa/2 + \gamma < 1/2$ .

For (C1), we note that had the  $\mathbf{x}_i$  been bounded by a constant, then  $E_{\boldsymbol{\beta}} \|\eta_{i,k}(\boldsymbol{\theta}, \boldsymbol{\beta})\|^2$  would have been bounded by a constant also. Since  $\max \|\mathbf{x}_i\|^2 = \mathcal{O}(p)$ , then  $E_{\boldsymbol{\beta}} \|\eta_{i,k}(\boldsymbol{\theta}, \boldsymbol{\beta})\|^2 =$

$\mathcal{O}(p)$  and  $E_{\boldsymbol{\beta}} \|\eta_i(\boldsymbol{\theta}, \boldsymbol{\beta})\|^2 = \mathcal{O}(Mp)$ , where  $Mp \leq cn^{\kappa+\gamma}$ . Therefore one can take  $n$  large enough such that  $C > \kappa + \gamma$  is satisfied with  $0 < d \leq 1$ . For (C2), we note that  $E \|\Psi_k(\boldsymbol{\beta}^*)\|^2 = \mathcal{O}(\max \|\mathbf{x}_i\|^2)$  and

$$\sum_{i=1}^n \sum_{k=1}^M E \|\Psi_k(\boldsymbol{\beta}^*)\|^2 = \mathcal{O}(Mnp) = \mathcal{O}(mn).$$

(C3) and (C4) are the hardest conditions to prove. As shown in Vanden Branden (2005), for  $\boldsymbol{\alpha} \in S_m$ ,

$$\boldsymbol{\alpha}^\top E [\Psi(\boldsymbol{\beta}) - \Psi(\boldsymbol{\beta}^*)] = n\boldsymbol{\alpha}^\top D_n(\boldsymbol{\beta} - \boldsymbol{\beta}^*) \quad (8)$$

$$+ \sum_{i=1}^n \sum_{k=1}^M \boldsymbol{\alpha}_k^\top \mathbf{x}_i \left\{ \tilde{f}'_{Y_i}(u)(\mathbf{x}_i^\top(\boldsymbol{\beta}(t_k) - \boldsymbol{\beta}^*(t_k)))^2 \right\} \quad (9)$$

$$+ \sum_{i=1}^n \sum_{k=1}^M \boldsymbol{\alpha}_k^\top \mathbf{x}_i \left\{ \sum_{l=1}^k d_{kl} \tilde{f}'_{C_i}(v)(\mathbf{x}_i^\top(\boldsymbol{\beta}(t_l) - \boldsymbol{\beta}^*(t_l)))^2 \right\} \quad (10)$$

where

$$d_{kl} = \begin{cases} -w_1 & l = 1 \\ w_{k-1} & l = k \\ -(w_l - w_{l-1}) & \text{otherwise} \end{cases}$$

$$d_{kli} = \begin{cases} d_{kk} \tilde{f}_{C_i}(\mathbf{x}_i^\top \boldsymbol{\beta}^*(t_k)) + \tilde{f}_{Y_i}(\mathbf{x}_i^\top \boldsymbol{\beta}^*(t_k)) & l = k \\ d_{kl} \tilde{f}_{C_i}(\mathbf{x}_i^\top \boldsymbol{\beta}^*(t_l)) & \text{otherwise} \end{cases}$$

and

$$nD_n = \begin{pmatrix} \sum_{i=1}^n d_{11i} \mathbf{x}_i \mathbf{x}_i^\top & \mathbf{0}_{p,p} & \dots & \dots & \dots & \dots & \mathbf{0}_{p,p} \\ \sum_{i=1}^n d_{21i} \mathbf{x}_i \mathbf{x}_i^\top & \sum_{i=1}^n d_{22i} \mathbf{x}_i \mathbf{x}_i^\top & \dots & \dots & \dots & \dots & \mathbf{0}_{p,p} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \sum_{i=1}^n d_{k1i} \mathbf{x}_i \mathbf{x}_i^\top & \sum_{i=1}^n d_{k2i} \mathbf{x}_i \mathbf{x}_i^\top & \dots & \sum_{i=1}^n d_{kki} \mathbf{x}_i \mathbf{x}_i^\top & \mathbf{0}_{p,p} & \dots & \mathbf{0}_{p,p} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \sum_{i=1}^n d_{M1i} \mathbf{x}_i \mathbf{x}_i^\top & \sum_{i=1}^n d_{M2i} \mathbf{x}_i \mathbf{x}_i^\top & \dots & \dots & \dots & \dots & \sum_{i=1}^n d_{MMi} \mathbf{x}_i \mathbf{x}_i^\top \end{pmatrix}. \quad (11)$$

Thus for (C3) to hold we require

$$\left| \sum_{i=1}^n \sum_{k=1}^M \boldsymbol{\alpha}_k^\top \mathbf{x}_i \left\{ \tilde{f}'_{Y_i}(u)(\mathbf{x}_i^\top(\boldsymbol{\beta}(t_k) - \boldsymbol{\beta}^*(t_k)))^2 + \sum_{l=1}^k d_{kl} \tilde{f}'_{C_i}(v)(\mathbf{x}_i^\top(\boldsymbol{\beta}(t_l) - \boldsymbol{\beta}^*(t_l)))^2 \right\} \right| = o(n^{1/2}) \quad (12)$$

or, as noted in Remark 2.3 of He and Shao (2000),

$$\begin{aligned} & \left| \sum_{i=1}^n \sum_{k=1}^M \boldsymbol{\alpha}_k^\top \mathbf{x}_i \left\{ \tilde{f}'_{Y_i}(u)(\mathbf{x}_i^\top(\boldsymbol{\beta}(t_k) - \boldsymbol{\beta}^*(t_k)))^2 + \sum_{l=1}^k d_{kl} \tilde{f}'_{C_i}(v)(\mathbf{x}_i^\top(\boldsymbol{\beta}(t_l) - \boldsymbol{\beta}^*(t_l)))^2 \right\} \right| \\ &= o((mn)^{1/2}). \end{aligned} \quad (13)$$

For (9) we have

$$\begin{aligned} & \sum_{i=1}^n \sum_{k=1}^M \boldsymbol{\alpha}_k^\top \mathbf{x}_i \tilde{f}'_{Y_i}(u)(\mathbf{x}_i^\top(\boldsymbol{\beta}(t_k) - \boldsymbol{\beta}^*(t_k)))^2 \\ & \leq \sum_{i=1}^n \sum_{k=1}^M |\boldsymbol{\alpha}_k^\top \mathbf{x}_i| \sum_{k=1}^M (\mathbf{x}_i^\top(\boldsymbol{\beta}(t_k) - \boldsymbol{\beta}^*(t_k)))^2 \\ & \leq \sum_{i=1}^n \|\mathbf{x}_i\| \left( \sum_{k=1}^M \|\boldsymbol{\alpha}_k\|^2 \right)^{1/2} \|\mathbf{x}_i\|^2 \sum_{k=1}^M \|\boldsymbol{\beta}(t_k) - \boldsymbol{\beta}^*(t_k)\|^2 \\ & = \mathcal{O}\left(\frac{m}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^3\right) = \mathcal{O}(p^{3/2}m) = \mathcal{O}(n^{5\gamma/2+\kappa}) \end{aligned}$$

and for (10)

$$\begin{aligned} & \sum_{i=1}^n \sum_{k=1}^M \boldsymbol{\alpha}_k^\top \mathbf{x}_i \sum_{l=1}^k d_{kl} \tilde{f}'_{C_i}(v)(\mathbf{x}_i^\top(\boldsymbol{\beta}(t_l) - \boldsymbol{\beta}^*(t_l)))^2 \\ & \leq \sum_{i=1}^n \sum_{k=1}^M \boldsymbol{\alpha}_k^\top \mathbf{x}_i d_{k1} \tilde{f}'_{C_i}(v)(\mathbf{x}_i^\top(\boldsymbol{\beta}(t_1) - \boldsymbol{\beta}^*(t_1)))^2 \end{aligned} \quad (14)$$

$$+ \sum_{i=1}^n \sum_{k=1}^M \boldsymbol{\alpha}_k^\top \mathbf{x}_i \sum_{l=2}^k d_{kl} \tilde{f}'_{C_i}(v)(\mathbf{x}_i^\top(\boldsymbol{\beta}(t_l) - \boldsymbol{\beta}^*(t_l)))^2 \quad (15)$$

Noting that  $d_{kl} = \mathcal{O}(1)$  for  $l = 1$  and  $\mathcal{O}(M)$  otherwise, we obtain

$$\begin{aligned} (14) & \leq \sum_{i=1}^n \left( \sum_{k=1}^M \|\boldsymbol{\alpha}_k^\top \mathbf{x}_i\|^2 \right)^{1/2} \left( \sum_{k=1}^M d_{k1}^2 \right)^{1/2} \|\mathbf{x}_i\|^2 \|\boldsymbol{\beta}(t_1) - \boldsymbol{\beta}^*(t_1)\|^2 \\ & \leq \sum_{i=1}^n \|\mathbf{x}_i\|^3 M^{1/2} \|\boldsymbol{\beta}(t_1) - \boldsymbol{\beta}^*(t_1)\|^2 = \mathcal{O}(p^{3/2} M^{1/2} \frac{m}{n}) \\ & = \mathcal{O}(n^{5\gamma/2+3\kappa/2}) \end{aligned}$$

and

$$\begin{aligned}
(15) &\leq \sum_{i=1}^n \sum_{k=1}^M \boldsymbol{\alpha}_k^\top \mathbf{x}_i \left( \sum_{l=2}^k d_{kl}^2 \right)^{1/2} \sum_{l=2}^k (\mathbf{x}_i^\top (\boldsymbol{\beta}(t_l) - \boldsymbol{\beta}^*(t_l)))^2 \\
&\leq \sum_{i=1}^n \sum_{k=1}^M \boldsymbol{\alpha}_k^\top \mathbf{x}_i \left( \frac{k-1}{M^2} \right)^{1/2} \|\mathbf{x}_i\|^2 \sum_{l=2}^k (\boldsymbol{\beta}(t_l) - \boldsymbol{\beta}^*(t_l))^2 \\
&= \mathcal{O}(M^{1/2} \frac{m}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^3) = \mathcal{O}(p^{3/2} M^{1/2} \frac{m}{n}) \\
&= \mathcal{O}(n^{5\gamma/2+3\kappa/2}).
\end{aligned}$$

With  $\gamma$  and  $\kappa$  satisfying  $2\gamma + \kappa < 1/2$ , the error from (C3) can be made  $o((mn)^{1/2}) = o(n^{\gamma/2+\kappa/2+1/2})$ .

(C4) needs to hold with  $A(n, m) = o(n/\log n)$ . The term  $\boldsymbol{\alpha}^\top \eta_i(\boldsymbol{\beta}, \boldsymbol{\beta}^*)$  is defined as

$$\boldsymbol{\alpha}^\top \eta_i(\boldsymbol{\beta}, \boldsymbol{\beta}^*) = \sum_{k=1}^M \boldsymbol{\alpha}_k^\top (\Psi_k(\mathbf{x}_i, \boldsymbol{\beta}) - \Psi_k(\mathbf{x}_i, \boldsymbol{\beta}^*) - E(\Psi_k(\mathbf{x}_i, \boldsymbol{\beta}) - \Psi_k(\mathbf{x}_i, \boldsymbol{\beta}^*))). \quad (16)$$

A Taylor series expansion for the expectation part of the expression gives

$$\sum_{k=1}^M \boldsymbol{\alpha}_k^\top \mathbf{x}_i \left[ \tilde{f}_{Y_i}(u)(\mathbf{x}_i^\top (\boldsymbol{\beta}(t_k) - \boldsymbol{\beta}^*(t_k))) + \sum_{l=1}^k d_{kl} \tilde{f}_{C_i}(v)(\mathbf{x}_i^\top (\boldsymbol{\beta}(t_l) - \boldsymbol{\beta}^*(t_l))) \right]$$

for some  $u$  and  $v$ . Similarly as for (C3) the first part of this term is bounded by  $\mathcal{O}((m/n)^{1/2} p) = \mathcal{O}(n^{\kappa/2+3\gamma/2-1/2})$  and the second part is bounded by  $\mathcal{O}(p(Mm/n)^{1/2}) = \mathcal{O}(n^{3\gamma/2+\kappa/2-1/2})$ .

Therefore

$$\boldsymbol{\alpha}^\top \eta_i(\boldsymbol{\beta}, \boldsymbol{\beta}^*) = \sum_{k=1}^M \boldsymbol{\alpha}_k^\top (\Psi_k(\mathbf{x}_i, \boldsymbol{\beta}) - \Psi_k(\mathbf{x}_i, \boldsymbol{\beta}^*)) + \mathcal{O}(n^{3\gamma/2+\kappa/2-1/2}).$$

This error term squared and multiplied by  $n$  is  $\mathcal{O}(n^{3\gamma+\kappa})$  which can be made  $o(n/\log n)$  if  $3\gamma+\kappa < 1$  so that it satisfies the requirement for (C4). For the term in  $\sum_{k=1}^M \boldsymbol{\alpha}_k^\top (\Psi_k(\mathbf{x}_i, \boldsymbol{\beta}) - \Psi_k(\mathbf{x}_i, \boldsymbol{\beta}^*))$  we introduce an indicator,  $I_{a_k, b_k}(Y)$ , with  $I_{a_k, b_k}(Y) = \pm 1$  if  $Y$  lies in between  $\mathbf{x}_i^\top a(t_k)$  and  $\mathbf{x}_i^\top b(t_k)$ , and 0 otherwise. Then

$$\begin{aligned}
\sum_{k=1}^M \boldsymbol{\alpha}_k^\top (\Psi_{ki}(\boldsymbol{\beta}) - \Psi_{ki}(\boldsymbol{\beta}^*)) &= \sum_{k=1}^M \boldsymbol{\alpha}_k^\top \mathbf{x}_i \text{sign}((\mathbf{x}_i^\top (\boldsymbol{\beta}(t_k) - \boldsymbol{\beta}^*(t_k))) \times \\
&[I(Y_i \leq C_i) I_{\boldsymbol{\beta}_k, \boldsymbol{\beta}_k^*}(Y_i) + I(Y_i > C_i) w(\boldsymbol{\beta}^*, t_k) I_{\boldsymbol{\beta}_k, \boldsymbol{\beta}_k^*}(C_i)] \\
&+ \sum_{k=1}^M \boldsymbol{\alpha}_k^\top \mathbf{x}_i I(Y_i > C_i) I(C_i \leq \mathbf{x}_i^\top \boldsymbol{\beta}(t_k)) (w_i(\boldsymbol{\beta}, t_k) - w_i(\boldsymbol{\beta}^*, t_k)).
\end{aligned}$$

The last term can be bounded using Lemma 2. For some constant  $D$

$$|w_i(\boldsymbol{\beta}, t_k) - w_i(\boldsymbol{\beta}^*, t_k)| = \left| \frac{(t_k - 1)(\tau_i(\boldsymbol{\beta}) - \tau_i(\boldsymbol{\beta}^*))}{(1 - \tau_i(\boldsymbol{\beta}))(1 - \tau_i(\boldsymbol{\beta}^*))} \right| \leq DT(n, m)$$

where  $T(n, m)$  is as defined in Lemma 2. Therefore the last term can be bounded by

$$\begin{aligned} \mathcal{O}(M^{1/2}p^{1/2}T(n, m)) &= \max(\mathcal{O}(M^{1/2}p(m/n)^{1/2}), \mathcal{O}(p^{1/2}/M^{1/2})) \\ &= \max(\mathcal{O}(n^{3\gamma/2+\kappa-1/2}), \mathcal{O}(n^{\gamma/2-\kappa/2})). \end{aligned}$$

Combining these results gives

$$\begin{aligned} |\boldsymbol{\alpha}^\top \boldsymbol{\eta}_i(\boldsymbol{\beta}, \boldsymbol{\beta}^*)| &\leq \sum_{k=1}^M \{|\boldsymbol{\alpha}_k^\top \mathbf{x}_i| [|I(Y_i \leq C_i)I_{\boldsymbol{\beta}_k, \boldsymbol{\beta}_k^*}(Y_i)| + |I(Y_i > C_i)I_{\boldsymbol{\beta}_k, \boldsymbol{\beta}_k^*}(C_i)|]\} \\ &\quad + \max(\mathcal{O}(n^{3\gamma/2+\kappa-1/2}), \mathcal{O}(n^{\gamma/2-\kappa/2})). \end{aligned}$$

This error term squared and multiplied by  $n$  will be  $o(n/\log n)$  if  $3\gamma + 2\kappa < 1$  and  $\gamma - \kappa < 0$ .

Finally for the term

$$\left( \sum_{k=1}^M \{|\boldsymbol{\alpha}_k^\top \mathbf{x}_i| [|I(Y_i \leq C_i)I_{\boldsymbol{\beta}_k, \boldsymbol{\beta}_k^*}(Y_i)| + |I(Y_i > C_i)I_{\boldsymbol{\beta}_k, \boldsymbol{\beta}_k^*}(C_i)|]\} \right)^2,$$

a bound is required on the number of observations for which  $I_{\boldsymbol{\beta}_k, \boldsymbol{\beta}_k^*}(Y_i)$  and  $I_{\boldsymbol{\beta}_l, \boldsymbol{\beta}_l^*}(Y_i)$  with  $l \neq k$  are both non-zero. By Lemma 2, this number is bounded by  $D^*T(n, m)M$  for some constant  $D^*$ . A bound of  $\mathcal{O}(p(m/n)^{1/2}) = \mathcal{O}(n^{3\gamma/2+\kappa/2-1/2})$  is thus obtained for the main part of the square. The cross term contributes

$$\mathcal{O}(p(m/n)^{1/2}T(n, m)M) = \max(\mathcal{O}(n^{5\gamma/2+5\kappa/2-1}), \mathcal{O}(n^{3\gamma/2+\kappa/2-1/2})).$$

The contribution of both terms can once again be made  $o(n/\log n)$  if  $5\gamma/2 + 5\kappa/2 < 1$  and  $3\gamma/2 + \kappa/2 < 1/2$ .

The constraints on  $\kappa$  and  $\gamma$  yield equations (4), (5) and (6).

All that is left is to verify that (C5) holds for these values.

According to Lemma 2.2 of He and Shao (2000), (C5) holds with the same  $A(n, m)$  as in (C4), provided that  $c_{n,m}^2 m \log n = \mathcal{O}(A(n, m))$ , where  $c_{n,m}$  is a sequence satisfying

$\sup_{\boldsymbol{\beta}, \mathbf{x}} \|\Psi(\mathbf{x}, \boldsymbol{\beta})\| \leq c_{n,m}$ . Here  $c_{n,m} = D^{**} M^{1/2} p^{1/2}$  for some constant  $D^{**}$ . Recalling that  $p = \mathcal{O}(n^\gamma)$ , it follows that  $c_{n,m}^2 m \log n = \mathcal{O}(A(n, m))$ , which concludes the proof of Theorem 3.1. ■

**Remark.** The results obtained in Theorem 3.1 are not optimal. For example, one possible choice for  $\gamma$  and  $\kappa$  is  $\gamma = 1/7$  and  $\kappa = 1/5$  which would give a rate of order  $n^{-23/35}$ . In addition, if condition (C4) holds with  $A(n, m) = o(\frac{n}{m \log n})$ , Theorem 2.2 of He and Shao (2000) gives asymptotic normality of the estimator, but requires tighter bounds than those obtained in Vanden Branden (2005), Neocleous et al (2006) and in Theorem 3.1. That is not to say that asymptotic normality is not possible. In fact, empirical results show that as the sample size  $n$  increases, the distribution of the CRQ-estimated  $\hat{\boldsymbol{\beta}}$  appears to approach a normal distribution.

Table 1: Comparison of performance for  $\beta_1(0.50)$  in the simulation model with approximate 20% censoring (Configuration 1). Knots at the quartiles of  $x_2$  were used for the spline terms.

$\tau = 0.50$	Bias	MAE	RMSE	ECP	EML
<b>n=200</b>					
lin	-0.00188	0.07646	0.11086	0.940	0.45406
pcs	-0.00012	0.00413	0.01115	0.996	0.04806
quad	0.00033	0.00436	0.00997	0.980	0.03552
cub	0.00024	0.00831	0.01420	0.968	0.05564
<b>n=500</b>					
lin	0.00262	0.05208	0.07554	0.936	0.28953
pcs	0.00003	0.00216	0.00419	0.990	0.01669
quad	-0.00019	0.00228	0.00452	0.950	0.01692
cub	-0.00003	0.00573	0.00843	0.960	0.03405
<b>n=1000</b>					
lin	-0.00198	0.03420	0.04850	0.952	0.20286
pcs	0.00001	0.00124	0.00228	0.982	0.00934
quad	-0.00011	0.00158	0.00291	0.950	0.01088
cub	-0.00005	0.00420	0.00609	0.954	0.02488

Table 2: Comparison of performance for  $\beta_1(0.75)$  in the simulation model with approximate 20% censoring (Configuration 1). Knots at the quartiles of  $x_2$  were used for the spline terms.

$\tau = 0.75$	Bias	MAE	RMSE	ECP	EML
<b>n=200</b>					
lin	-0.00167	0.06349	0.10313	0.928	0.40821
pcs	0.00081	0.00784	0.01576	0.969	0.05667
quad	-0.00004	0.00332	0.00787	0.994	0.03060
cub	0.00033	0.00637	0.01171	0.969	0.05071
<b>n=500</b>					
lin	0.00349	0.04290	0.06439	0.940	0.25481
pcs	-0.00001	0.00436	0.00771	0.949	0.02945
quad	-0.00014	0.00169	0.00352	0.978	0.01349
cub	-0.00028	0.00411	0.00707	0.966	0.02916
<b>n=1000</b>					
lin	-0.00432	0.03272	0.04355	0.954	0.17951
pcs	-0.00017	0.00353	0.00508	0.946	0.02003
quad	-0.00005	0.00124	0.00209	0.964	0.00815
cub	-0.00002	0.00302	0.00480	0.968	0.01980

Table 3: Comparison of performance for  $\beta_1(0.50)$  in the simulation model with  $n = 500$  and approximate 40% censoring (Configuration 2). Knots at (a) the 33th and 66th quantiles, (b) the quartiles and (c) the 20th, 40th, 60th and 80th quantiles of  $x_2$  were used for the spline terms.

$\tau = 0.50$	Bias	MAE	RMSE	ECP	EML
<b>Linear term in <math>x_2</math></b>	-0.1074	0.1069	0.1256	0.6640	0.2835
<b>Piecewise linear spline</b>					
(a)	-0.0166	0.0173	0.0233	0.7980	0.0645
(b)	0.0108	0.0109	0.0212	0.9457	0.0641
(c)	0.0056	0.0081	0.0144	0.9618	0.0526
<b>Quadratic spline</b>					
(a)	0.0276	0.0288	0.0348	0.7560	0.0917
(b)	0.0010	0.0032	0.0055	0.9739	0.0219
(c)	0.0030	0.0047	0.0081	0.9379	0.0279
<b>Cubic spline</b>					
(a)	0.0018	0.0038	0.0060	0.9700	0.0242
(b)	0.0061	0.0080	0.0110	0.9280	0.0379
(c)	0.0008	0.0026	0.0040	0.9699	0.0172

Table 4: Comparison of performance for  $\beta_1(0.75)$  in the simulation model with  $n = 500$  and approximate 40% censoring (Configuration 2). Knots at (a) the 33th and 66th quantiles, (b) the quartiles and (c) the 20th, 40th, 60th and 80th quantiles of  $x_2$  were used for the spline terms.

$\tau = 0.75$	Bias	MAE	RMSE	ECP	EML
<b>Linear term in <math>x_2</math></b>	-0.2084	0.2116	0.2239	0.3260	0.3246
<b>Piecewise linear spline</b>					
(a)	-0.0247	0.0253	0.0330	0.7818	0.0918
(b)	-0.0033	0.0091	0.0135	0.9277	0.0491
(c)	0.0023	0.0052	0.0093	0.9351	0.0361
<b>Quadratic spline</b>					
(a)	0.0111	0.0104	0.0159	0.8741	0.0500
(b)	0.0011	0.0033	0.0050	0.9834	0.0210
(c)	0.0021	0.0048	0.0077	0.9436	0.0289
<b>Cubic spline</b>					
(a)	0.0013	0.0040	0.0063	0.9529	0.0246
(b)	0.0035	0.0052	0.0081	0.9306	0.0306
(c)	0.0011	0.0029	0.0044	0.9741	0.0176

Table 5: Comparison of performance for  $Q(\tau | x)$  in the simulation model with  $n = 500$  and approximate 40% censoring (Configuration 2). Knots at (a) the 33th and 66th quantiles, (b) the quartiles and (c) the 20th, 40th, 60th and 80th quantiles of  $x_2$  were used for the spline terms.

	$\tau = 0.50$		$\tau = 0.75$	
	RMSE	Bias	RMSE	Bias
<b>Linear term in <math>x_2</math></b>	1.3968	0.3233	1.3275	-0.0271
<b>Piecewise linear spline</b>				
(a)	0.4967	0.2426	0.4875	-0.1078
(b)	0.6699	0.3322	0.6281	-0.0182
(c)	0.4775	0.2711	0.4842	-0.0793
<b>Quadratic spline</b>				
(a)	0.8891	0.4601	0.7975	0.1096
(b)	0.4731	0.2755	0.4721	-0.0750
(c)	0.5090	0.2804	0.4973	-0.0700
<b>Cubic spline</b>				
(a)	0.4784	0.2682	0.4766	-0.0822
(b)	0.6032	0.3072	0.5558	-0.0432
(c)	0.4956	0.2680	0.4991	-0.0824

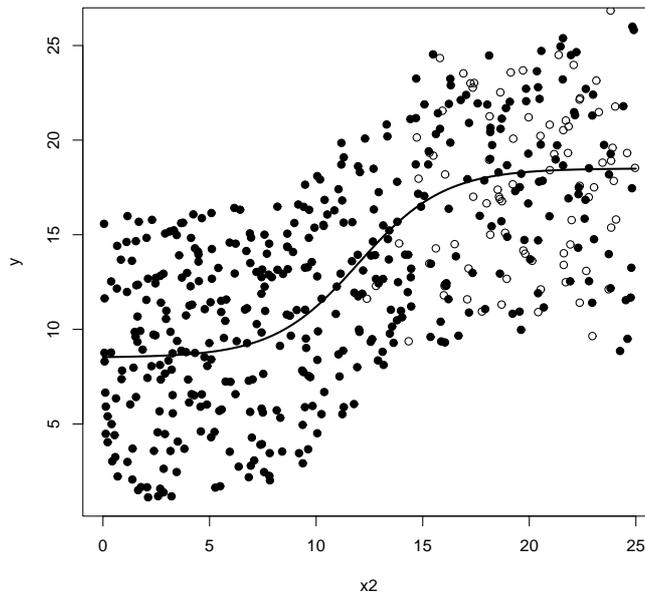


Figure 1: Scatterplot of Configuration 1 used in the simulation experiment. Censored points are shown as open circles, uncensored points as filled circles. The conditional median line evaluated at the mean of  $x_1$  is also shown.

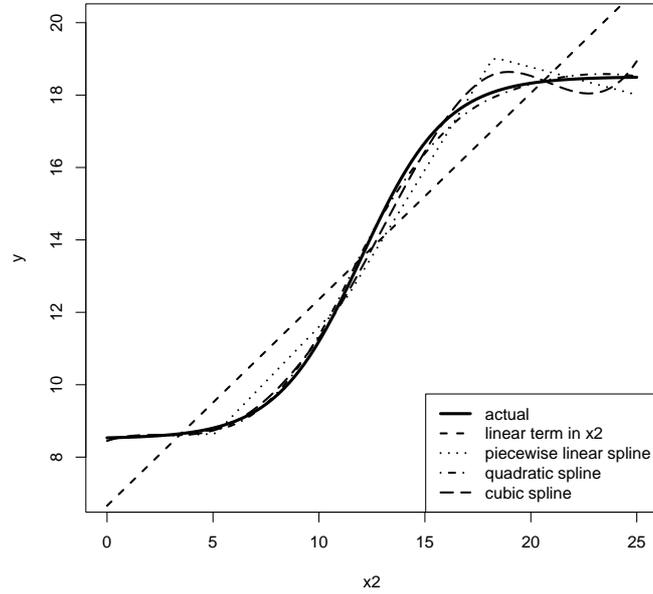


Figure 2: Various model fits for the nonlinear term in the simulation experiment (Configuration 1). Shown here are the actual median (solid line) and model-estimated conditional median lines (dashed or dotted) evaluated at the mean of  $x_1$ .

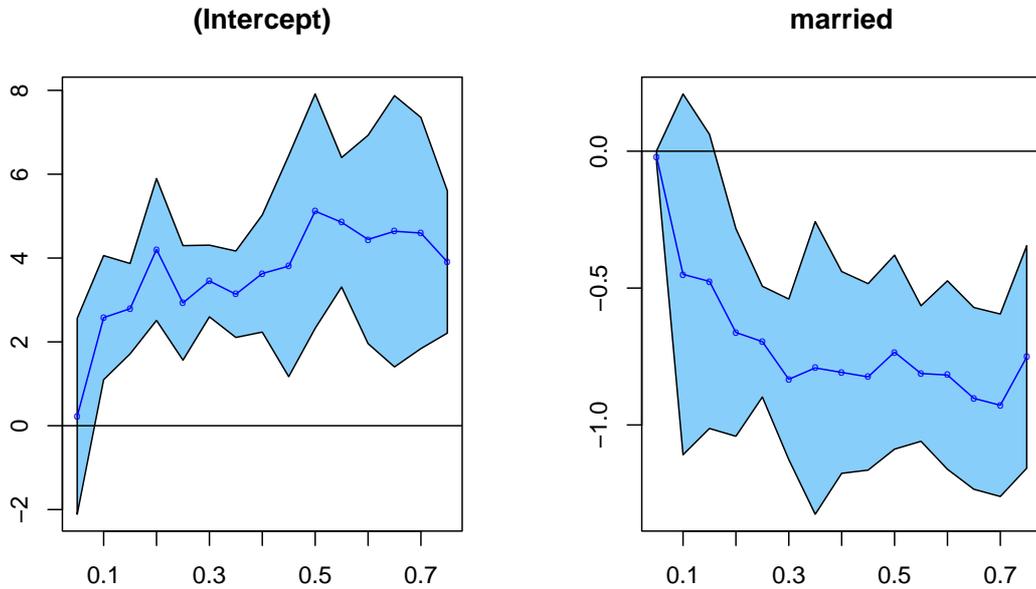


Figure 3: Estimated linear coefficients  $\hat{\beta}_0(\tau)$  and  $\hat{\beta}_1(\tau)$  in model (7) with 95% bootstrap pointwise confidence intervals plotted against  $\tau$  for  $0 < \tau \leq 0.75$ .

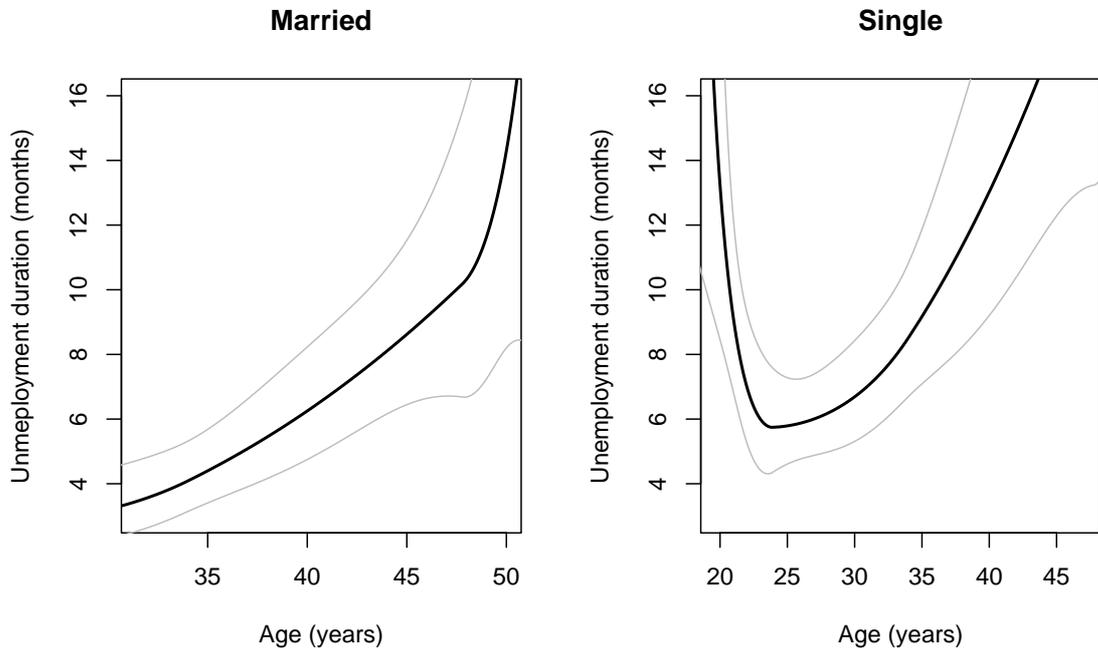


Figure 4: Estimated median unemployment duration against age for German males. The black line shows the median, grey lines show 95% pointwise confidence limits.

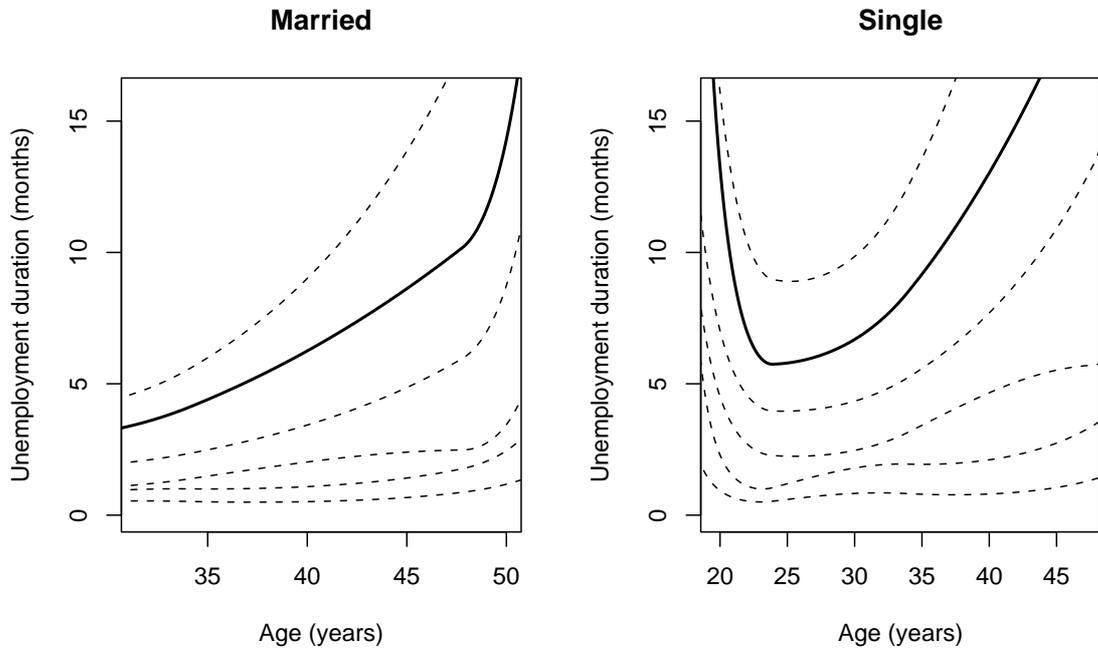
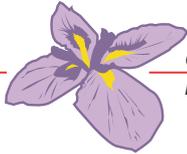


Figure 5: Estimated deciles of unemployment duration against age for German males. The solid line shows the median, dashed lines show the other deciles from 1st to 6th.



## IRISS Working Papers

The IRISS Working Paper Series has been created in 1999 to ensure a timely dissemination of the research outcome from the IRISS-C/I programme. They are meant to stimulate discussion and feedback. The working papers are contributed by CEPS/INSTEAD resident staff, research associates and visiting researchers.

### The fifteen most recent papers

Neocleous T. & Portnoy S., 'A Partially Linear Censored Quantile Regression Model for Unemployment Duration', IRISS WP 2008-07, September 2008.

Kankara M. & Moors G., 'Measurement Equivalence and Extreme Response Bias in the Comparison of Attitudes across Europe', IRISS WP 2008-06, May 2008.

Prejmerean M. & Vasilache S., 'What's a university worth? Changes in the lifestyle and status of post-2000 European Graduates.', IRISS WP 2008-05, February 2008.

Takhtamanova Y. & Sierminska E., 'Gender differences in the effect of monetary policy on employment: The case of nine OECD countries.', IRISS WP 2008-04, February 2008.

Tamilina L., 'The analysis of welfare state effects on social trust in a multidimensional approach', IRISS WP 2008-03, February 2008.

Corsini L., 'Institutions, Technological Change and the Wage Differentials Between Skilled and Unskilled Workers: Theory and Evidence from Europe', IRISS WP 2008-02, January 2008.

Gerber P. & Fleuret S., 'Cartographier une enquête à l'échelle intra-urbaine: bien-être et personnes âgées de la ville de Luxembourg', IRISS WP 2008-01, January 2008.

Pavlopoulos D., Fouarge D., Muffels R. & Vermunt J., 'Who benefits from a job change: The dwarfs or the giants?', IRISS WP 2007-16, December 2007.

Martin L., 'The impact of technological changes on incentives and motivations to work hard', IRISS WP 2007-15, December 2007.

Popescu L., Rat C. & Rebeleanu-Bereczki A., 'Self-Assessed Health Status and Satisfaction with Health Care Services in the Context of the Enlarged European Union', IRISS WP 2007-14, November 2007.

Weziak D., 'Measurement of national intellectual capital application to EU countries', IRISS WP 2007-13, November 2007.

D'Angelo E. & Lilla M., 'Is there more than one linkage between Social Network and Inequality?', IRISS WP 2007-12, November 2007.

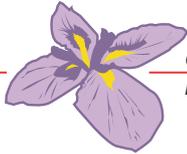
Lilla M., 'Income Inequality and Education Premia', IRISS WP 2007-11, November 2007.

Stanciole A., 'Health Insurance and Life Style Choices: Identifying the Ex Ante Moral Hazard', IRISS WP 2007-10, November 2007.

Raileanu Szeles M., 'The patterns and causes of social exclusion in Luxembourg', IRISS WP 2007-09, August 2007.

### Electronic versions

Electronic versions of all IRISS Working Papers are available for download at  
<http://www.ceps.lu/iriss/wps.cfm>



*IRISS-C/I is a visiting researchers programme at CEPS/INSTEAD, a socio-economic policy and research centre based in Luxembourg. It finances and organises short visits of researchers willing to undertake empirical research in economics and other social sciences using the archive of micro-data available at the Centre.*

### What is offered?

In 1998, CEPS/INSTEAD has been identified by the European Commission as one of the few *Large Scale Facilities* in the social sciences, and, since then, offers researchers (both junior and senior) the opportunity to spend time carrying out their own research using the local research facilities. This programme is currently sponsored by the European Community's 6th Framework Programme. Grants cover travel expenses and on-site accommodation. The expected duration of visits is in the range of 2 to 12 weeks.

### Topics

The major resource offered to visitors is access to a series of internationally comparable longitudinal surveys on living conditions at the household and individual level. The anonymised micro-data provide information on wages and income, health, education, employment and professional activities, accommodation, social relations,... Comparable micro-data are available for EU countries, Central European countries, as well as the USA. These data offer opportunities to carry out research in fields such as *survey and panel data methodology, income distribution and welfare, income and poverty dynamics, multi-dimensional indicators of poverty and deprivation, gender, ethnic and social inequality, unemployment and labour supply behaviour, education and training, social protection and redistributive policies, fertility and family structures, new information technologies in households and firms, ...*

### Who may apply?

All individuals (doctoral students as well as experienced academics) conducting research in an institution within the EU-25 or an FP6 Associated State. IRISS-C/I can be meeting place for groups of researchers working on a joint project. We therefore encourage joint proposals by two or more researchers.

For more detailed information and application form, please consult our website: <http://www.ceps.lu/iriss> or contact us at

IRISS-C/I, CEPS/INSTEAD  
BP 48, L-4501 Differdange, G.-D. Luxembourg  
Tel: +352 585855 610; Fax: +352 585588  
E-mail: [iriss@ceps.lu](mailto:iriss@ceps.lu)