

Fast Inference for Quantile Regression with Tens of Millions of Observations

Sokbae Lee, Yuan Liao, Myung Hwan Seo & Youngki Shin
Columbia Univ. Rutgers Univ. Seoul National Univ. McMaster Univ.

University of California, Riverside
December 6, 2022

Motivation

- We consider a quantile regression, cross-sectional data.

$$y_i = x_i' \beta + e_i, \quad P(e_i \leq 0 | x_i) = \tau$$

- with more emphasis on inference.
- **This paper:** $n \sim 10^7$.
- Standard M-estimation would involve:
 - (1) solving optimization problems
 - (2) estimate sandwich matrices
- Neither would be an easily scalable task for datasets of such sizes.

Motivation

- IPUMS USA dataset collects U.S. census microdata from 2000 to 2020. 3.5 million households/year
- Most studies use much smaller subsamples, with two recommended features by IPUMS:
 - (1) “Select cases”: only particular kinds of households
 - (2) “Customize sample sizes” : draw smaller random subsets.
- Researchers also manually filter the information for sample selections.

Motivation

- All these may introduce sample selection biases.
- IPUMS reminded cautions to researchers throughout the users' selection sessions.
 - *"the population at risk for answering the question – can differ subtly or markedly across samples."*
 - *"Users should be careful with the case selection feature... thereby inadvertently excluding those samples from your dataset."*

This paper

- Propose a stochastic (sub-) gradient descent (SGD) based method, with advantages of:
 - **fast**: less than 10 seconds on a regular PC for $n \sim 10^7$.
 - **memory-efficient**: computed recursively, do not store many observations in the computation.
- How it works:
 - SGD update gives a solution path:

$$\underbrace{\beta_0}_{\text{initial}}, \quad \beta_1, \dots, \beta_n$$

$$\text{Polyak-Ruppert average: } \hat{\beta} := \frac{1}{n} \sum_{i=1}^n \beta_i.$$

- Standardize the estimator by a transformation of the cumulative sums, denoted by \hat{V}_n , which renders a pivotal statistic.

Gradient Descent

Let β^* be the parameter of interest:

$$\beta^* := \arg \min_{\beta \in \mathbb{R}^d} Q(\beta)$$

where $Q := E[q(\beta, Y)]$ and q is diff. and convex. Let $\{Y_t\}_{t=1}^n$ be a random sample. The sample analogue of the FOC is

$$\frac{1}{n} \sum_{t=1}^n \nabla q(\hat{\beta}, Y_t) = 0.$$

If we cannot solve it directly, the solution is computed iteratively:

$$\hat{\beta}_m = \hat{\beta}_{m-1} - \gamma_m \frac{1}{n} \sum_{t=1}^n \nabla q(\hat{\beta}_{m-1}, Y_t).$$

Stochastic Gradient Descent

Limitations of gradient descent:

- It calculates the derivatives for the entire dataset.
- It requires a larger memory size as the dataset increases.

Binding time budget or the memory size occurs more often in modern empirical applications.

Robbins and Monro (1951) proposed the stochastic gradient descent (SGD) solution path as

$$\beta_t = \beta_{t-1} - \gamma_t \nabla q(\beta_{t-1}, Y_t).$$

SGD has advantages when we face a [large-scale dataset](#) or [online machine learning](#).

Examples: Chen and White (2002), Khan, Lan, and Tamer (2021).

SGD Averaging

Recall that we aim to develop **online inference** with SGD.

We study the classical Polyak-Ruppert averaging estimator (Polyak (1990) and Ruppert(1988)): $\bar{\beta}_n := n^{-1} \sum_{t=1}^n \beta_t$.

Polyak and Juditsky (1992) established regularity conditions under which the averaging estimator $\bar{\beta}_n$ is asymptotically normal:

$$\sqrt{n} (\bar{\beta}_n - \beta^*) \xrightarrow{d} \mathcal{N}(0, \Upsilon),$$

where the asymptotic variance Υ has a sandwich form

$$\Upsilon := H^{-1} S H^{-1},$$

and $H := \nabla^2 Q(\beta^*)$ is the Hessian matrix and $S := \mathbb{E} [\nabla q(\beta^*, Y) \nabla q(\beta^*, Y)']$ is the score variance.

SGD Averaging in Online Learning

In online learning, data arrive sequentially.

The Polyak-Ruppert estimator $\bar{\beta}_n$ can be computed recursively by the updating rule

$$\bar{\beta}_t = \bar{\beta}_{t-1} \frac{t-1}{t} + \frac{\beta_t}{t},$$

which implies that it is well suited to the online setting.

Examples include

- Linear regression (with a large dataset)
- Logistic regression
- Quantile regression (using a subgradient)

Linear Quantile Regression I

- Consider the linear quantile regression

$$y_t = x_t' \beta_u^* + \varepsilon_t, \quad P(\varepsilon_t \leq 0 | x_t) = u \in (0, 1).$$

- The population loss is $Q(\beta) = \mathbb{E}[q(\beta; x_t, y_t)]$, with the check function

$$q(\beta; x, y) = (y - x' \beta)(u - 1\{y - x' \beta \leq 0\}).$$

- As the check function q is not differentiable, we propose to use its subgradient:

$$\nabla q(\beta; x, y) = -x[u - 1\{y \leq x' \beta\}],$$

leading to the online update:

$$\beta_t = \beta_{t-1} - \gamma_t \nabla q(\beta_{t-1}, x_t, y_t).$$

Linear Quantile Regression II

- $\bar{\beta}_n$ is asymptotically normal with asymptotic variance $\Upsilon = H^{-1}SH^{-1}$, where

$$S = \mathbb{E}x_t x_t' u(1-u) \quad \text{and} \quad H = \mathbb{E}x_t x_t' f_\varepsilon(0|x_t),$$

where $f_\varepsilon(\cdot|x_t)$ denotes the conditional density of ε_t given x_t .

- History: Koenker and Bassett (1978)'s seminal work and application to wage regression by Chamberlain (1994), Buchinsky (1994, 1998), Angrist et al. (2006), etc
- Computation: Portnoy and Koenker (1997)'s linear programming through interior-point algorithms and convolution type smoothing by Fernandes et al. (2021), He et al. (2021), Tan et al. (2022).
- He et al.'s **conquer** has received attention and boast its capability to make inference with $(n, p) = (4000, 100)$.

Inference Methods

Online Inference

Although the asymptotic normality result by Polyak and Juditsky (1992) was established about three decades ago, it is only past few years that online inference has gained increasing interest in the literature.

It is challenging to estimate the asymptotic variance Υ in an online fashion.

This is because the naive implementation of estimating it requires storing all data, thereby losing the advantage of online learning.

Method 1: Plug-In

Chen et al. (2020) addressed this issue by estimating H and S using the online iterated estimator β_t , and recursively updating them whenever a new observation is available.

However, the plug-in estimator requires that the Hessian matrix be computed to estimate H .

In other words, it is necessary to have strictly more inputs than the SGD solution paths β_t to carry out inference. It is the case even when a t-statistic is computed for each regression coefficient.

In applications, it can be demanding to compute the Hessian matrix and its inverse.

They do not cover the quantile regression.

Method 2: Batch-Means

This method proposed by Chen et al. (2020) and Zhu et al. (2021) directly estimates the variance of the averaged online estimator $\bar{\beta}_n$ by dividing $\{\beta_1, \dots, \beta_n\}$ into batches with increasing batch size.

The batch-means estimator is based on the idea that correlations among batches that are far apart decay exponentially fast; therefore, one can use nonparametric empirical covariance to estimate Υ .

However, this approach requires the batch size should increase exponentially fast, and it turns out that the performance is not satisfactory.

Weakly Dependent Batches

- Zhu et al (2021) take batches: e.g.,

$$\begin{aligned} B_1 &= \{\beta_1\}, & B_2 &= \{\beta_1, \beta_2\}, \\ B_3 &= \{\beta_3\}, & B_4 &= \{\beta_3, \beta_4\}, \\ B_5 &= \{\beta_5\}, & B_6 &= \{\beta_5, \beta_6\}, \end{aligned}$$

- In fact, the gap $B_1, B_3, B_5 \dots$ should be much larger, so that they are weakly dependent. So the “batch gap” and “batch size” are tuning parameters.
- The “weakly dependent” batches are required, and the goal is the average of blockwise sample variance

$$\hat{\Upsilon} \rightarrow^P \Upsilon.$$

- Then

$$\sqrt{n} \hat{\Upsilon}^{-1/2} (\bar{\beta}_n - \beta^*) \rightarrow^d N(0, I).$$

Method 3: Bootstrap

Instead of estimating the asymptotic variance, Fang et al. (2018) proposed a bootstrap procedure for online inference.

Specifically, they proposed to use a large number (say, B) of randomly perturbed SGD solution paths: for all $b = 1, \dots, B$, starting with $\beta_0^{(b)} = \beta_0$ and then iterating

$$\beta_t^{(b)} = \beta_{t-1}^{(b)} - \gamma_t \eta_t^{(b)} \nabla q \left(\beta_{t-1}^{(b)}, Y_t \right),$$

where $\eta_t^{(b)} > 0$ is an independent and identically distributed random variable that has mean one and variance one.

The bootstrap procedure needs strictly more inputs than computing $\bar{\beta}_n$ and can be time-consuming.

Our Approach: Random Scaling

- Lee, Liao, Seo, Shin (LLSS, 2022) proposed not to estimate the asymptotic variance Υ , but to studentize $\sqrt{n}(\bar{\beta}_n - \beta^*)$, or its each element, via $\hat{V}_n^{1/2}$, where

$$\hat{V}_n := \frac{1}{n} \sum_{s=1}^n \left\{ \frac{1}{\sqrt{n}} \sum_{t=1}^s (\beta_t - \bar{\beta}_n) \right\} \left\{ \frac{1}{\sqrt{n}} \sum_{t=1}^s (\beta_t - \bar{\beta}_n) \right\}',$$

or by its corresponding diagonal term, $\hat{V}_{n,jj}^{1/2}$.

- It converges in distribution to a pivotal distribution up to an unknown scale, which is the same as the asymptotic variance of the average SGD estimator. This leverages insights from the time series literature (e.g. Kiefer et al. (2000)).
- It has been already being adopted in other machine-learning literature like in federated learning (Li et al. 2022), Kiefer-Wolfowitz method (Chen et al. 2021) among others.

Algorithm

Input: function $\nabla q(\cdot)$, parameters (γ_0, a) for step size $\gamma_t = \gamma_0 t^{-a}$ for $t \geq 1$

Initialize: set initial values for $\beta_0, \bar{\beta}_0, A_0, b_0$

for $t = 1, 2, \dots$ **do**

Receive: new observation Y_t

$$\beta_t = \beta_{t-1} - \gamma_t \nabla q(\beta_{t-1}, Y_t)$$

$$\bar{\beta}_t = \bar{\beta}_{t-1} \frac{t-1}{t} + \frac{\beta_t}{t}$$

$$A_t = A_{t-1} + t^2 \bar{\beta}_t \bar{\beta}_t'$$

$$b_t = b_{t-1} + t^2 \bar{\beta}_t$$

$$c_t = c_{t-1} + t^2$$

 Obtain \hat{V}_t by

$$\hat{V}_t = t^{-2} (A_t - \bar{\beta}_t b_t' - b_t \bar{\beta}_t' + \bar{\beta}_t \bar{\beta}_t' c_t)$$

Output: $\bar{\beta}_t, \hat{V}_t$

end

Criteria for Online Inference Methods

Table: Criteria for Online Inference Methods

Method	FXY (18) Bootstrap	CLTZ (20) Plug-In	CLTZ (20) Batch Means	ZCW (21) Batch Means	LLSS (22) Random Scaling
Is it possible					
to avoid resampling?		✓	✓	✓	✓
to avoid Hessian?	✓		✓	✓	✓
to update recursively?	✓	✓		✓	✓

Note. FXY (18), CLTZ (20), and ZCW (21) refer to Fang et al. (2018), Chen et al. (2020), and Zhu et al. (2021), respectively.

Tuning parameters

- Initialize from smoothed QR from Kaplan and Sun (2017), Fernandes et al (2021), He et al (2021).
- specifically, computed fast using gradient descent update; R package: conquer and 5% of the data to compute β_0 ,
- stepsize γ_t is set closer to upper bound, $t^{-0.501}$.

Inference for Sub-vectors

- Empirical studies often involve many controls: β^* is “long”.
- But of interest is a sub-vector, containing only 1~2 elements.
- It is straightforward to cast sub-vector inference:

$$\beta_i = \text{full vector} \quad (1)$$

$$\bar{\beta}_i = \text{sub vector} \quad (2)$$

$$\widehat{V}_i = \text{sub matrix} \quad (3)$$

Forneron and Ng (2020)

- They proposed

$$\hat{\beta}_i = \hat{\beta}_{i-1} - \gamma H^{-1} \nabla q(Y_i, \hat{\beta}_{i-1})$$

- Can consistently estimate the variance using

$$\frac{\gamma^2}{1 - (1 - \gamma)^2} \frac{1}{n} \sum_i (\hat{\beta}_i - \bar{\beta}_n)(\hat{\beta}_i - \bar{\beta}_n)'$$

where $\bar{\beta}_n = \frac{1}{n} \sum_i \hat{\beta}_i$.

- The use of a fixed step size and the Hessian H^{-1} is the key for the consistency:

$$\hat{\beta}_i \approx \rho \hat{\beta}_{i-1} + \text{noise}$$

- Representable as a stationary AR(1) model.

Extensions to cluster-dependent QR

Consider

$$y_{i,t} = x'_{i,t}\beta + \epsilon_{i,t}$$

where $i = 1 \dots n$ are “clusters”, and $t = 1, \dots, T_i$ are individuals within clusters.

- Allow arbitrary dependence within cluster; Independence between clusters.
- Large n , bounded T_i .
- current theory requires $T_i = T$ for all i .
- A generalization in the spirit of Hansen and Lee (2019, JoE) would be very interesting.

Extensions to cluster-dependent QR

- Hagemann 2017 proposed a multiplier bootstrap for this setting.
- Our proposal: recursively update

$$\beta_i = \beta_{i-1} - \gamma_i \frac{1}{T_i} \sum_{t=1}^{T_i} \nabla q(Y_{i,t}, \beta_i)$$

where $\nabla q(Y, \beta) = -x(\tau - 1\{y < x'\beta\})$.

- Final estimator: $\bar{\beta}_n$
- Inference: Same random-scaling \hat{V}_n .

A random clustering ?

- 1 Consider a random sampling from $W = (D_1, D_2, X_1, X_2, \dots)$.
- 2 then sampling scheme for a cluster such that each cluster $(X_{g1}, \dots, X_{g,n_g})$ for $g = 1, 2, \dots, G$ is a transformation such that

$$(X_{g1}, \dots, X_{g,n_g}) = (X_{D_1}, \dots, X_{D_2}),$$

where $n_g = D_2 - D_1 + 1$.

- 3 Here, we assume $D_1 \leq D_2$ but we do not assume much about other dependence on the vector W . Then,

$$\bar{X}_g = n_g^{-1} \sum_{i=1}^{n_g} X_{g,i}$$

is an iid sequence. That is, \bar{X}_g is iid over $g = 1, 2, \dots$ unconditionally, while it is not conditionally on $(D_{g,1}, D_{g,2})$.

Theoretical Results

Functional Central Limit Theorem for Online SGD

We first extend Polyak and Juditsky (1992)'s central limit theorem (CLT) to a *functional* CLT (FCLT) for partial sum process:

$$\frac{1}{\sqrt{n}} \sum_{t=1}^{[nr]} (\beta_t - \beta^*) \Rightarrow \Upsilon^{1/2} W(r), \quad r \in [0, 1],$$

where \Rightarrow stands for the weak convergence in $\ell^\infty[0, 1]$ and $W(r)$ stands for a vector of the independent standard Wiener processes on $[0, 1]$.

That is, the partial sum of the online updated estimates β_t converges weakly to a rescaled Wiener process,

Note that the scaling $\Upsilon^{1/2}$ is equal to the square root asymptotic variance of the Polyak-Ruppert average.

PJ's approximation

Their CLT is built on a brilliant stochastic approximation that

$$\frac{1}{\sqrt{n}} \sum_{t=1}^{[n]} (\beta_t - \beta^*) = \frac{1}{\sqrt{n}} \sum_{t=1}^{[n]} H^{-1} \xi_t + o_p(1),$$

where ξ_t is an mds sequence whose variance converges to the score variance.

For FCLT, it is required to extend the approximation to the uniform approximation.

Conditions for Quantile Regression

Let the partial derivative $\frac{d}{de} f_\epsilon(\cdot|x_i)$ exist and assume that

- i) there exist positive constants ϵ and c_0 such that

$$\inf_{|\beta - \beta^*| < \epsilon} \lambda_{\min} \left(\mathbb{E}[x_i x_i' f_\epsilon(x_i'(\beta - \beta^*)|x_i)] \right) > c_0,$$

- ii) $\sup_b \mathbb{E}[\|x_i\|^3 A(b, x_i)] < C$ for some constant $C < \infty$, where $A(b, x_i) := \left| \frac{d}{de} f_\epsilon(x_i' b | x_i) \right| + f_\epsilon(x_i' b | x_i)$,
- iii) $\mathbb{E}[(\|x_i\|^6 + 1) \exp(\|x_i\|^2)] < C$ for some constant $C < \infty$,

This is a set of low-level conditions to meet Gadat and Panloup (2022)'s consistency without strong convexity and an extension of LLSS (2022)'s FCLT.

Main Theorem

Let for any $\ell \leq d$ linear restrictions

$$H_0 : R\beta^* = c,$$

where R is an $(\ell \times d)$ -dimensional known matrix of rank ℓ and c is an ℓ -dimensional known vector.

Theorem

Suppose $\text{rank}(R) = \ell$. Under the stated Assumptions and H_0 ,

$$\begin{aligned} & n (R\bar{\beta}_n - c)' \left(R\hat{V}_n R' \right)^{-1} (R\bar{\beta}_n - c) \\ & \xrightarrow{d} W(1)' \left(\int_0^1 \bar{W}(r) \bar{W}(r)' dr \right)^{-1} W(1), \end{aligned}$$

where $W(\cdot)$ is an ℓ -dimensional vector of the standard Wiener processes and $\bar{W}(r) := W(r) - rW(1)$.

Special Case: t-Statistic

the t-statistic for each j converges in distribution:

$$\frac{\sqrt{n} \left(\bar{\beta}_{n,j} - \beta_j^* \right)}{\sqrt{\hat{V}_{n,jj}}} \xrightarrow{d} W_1(1) \left[\int_0^1 \{W_1(r) - rW_1(1)\}^2 dr \right]^{-1/2},$$

- The asymptotic distribution is mixed normal and symmetric around zero,
- It is the same as the distribution of the statistics observed in the estimation of the cointegration vector by Johansen (1991). They are different statistics but have the identical distribution as functions of the standard Wiener process as shown by Abadir and Paruolo (2002).
- Abadir and Paruolo (1997) obtained a closed form density function.

Monte Carlo Experiments

MC Simulations

- Setting: $\dim(x_t) = d \in \{10, 30, 180, 320, 1000\}$,
 $n \in \{10^5, 10^6, 10^7\}$.
- Compare 5 methods:
 - Proposed** : proposed
 - QR** : “standard”, R package quantreg.
 - Conquer-plugin** : Fernandes et al (2021), R package conquer.
 - Conquer-bootstrap** : He et al (2021), R package conquer.
 - SGD-bootstrap** : Fang et al et al (2016). Bootstrap-online learning
- Constraints: 10 hour and 192 Gb RAM for single replication

Can you compute?

Proposed : no pressure

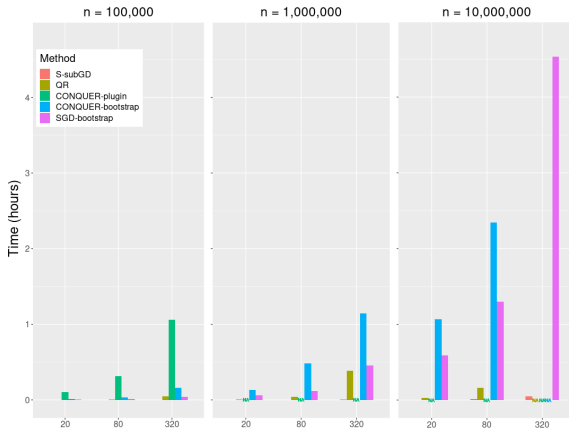
QR : out-of memory when $d = 320$ and $n \sim 10^7$

Conquer-plugin : out-of memory when $n \sim 10^6$

Conquer-bootstrap : out-of-time when $d = 320$ and $n \sim 10^7$

SGD-bootstrap : barely survived

Figure: Computation time



Note: Observe 'NA' for several cases.

Figure: Computation time

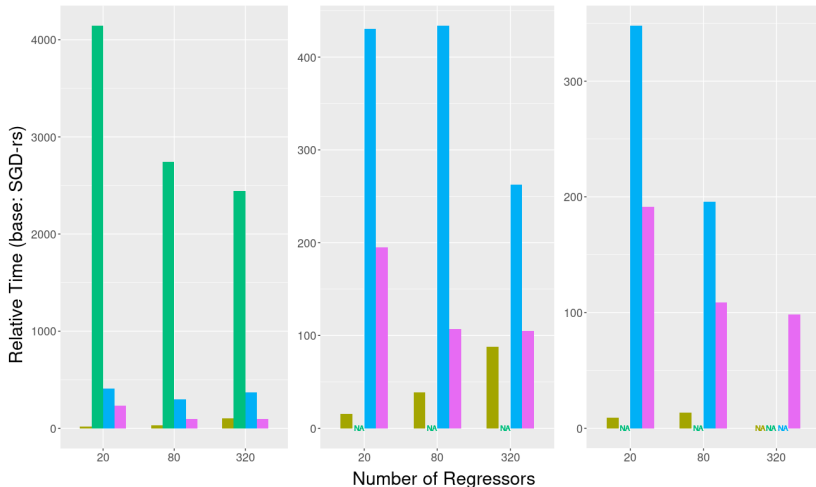


Figure: Coverage Rate

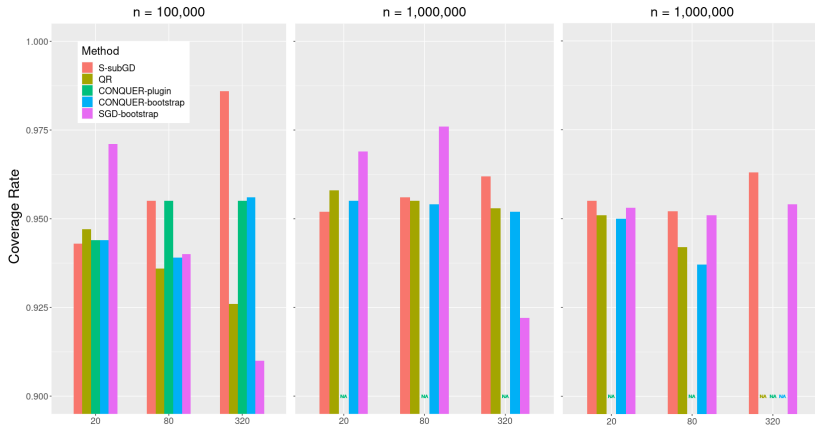


Figure: Confidence Interval Length

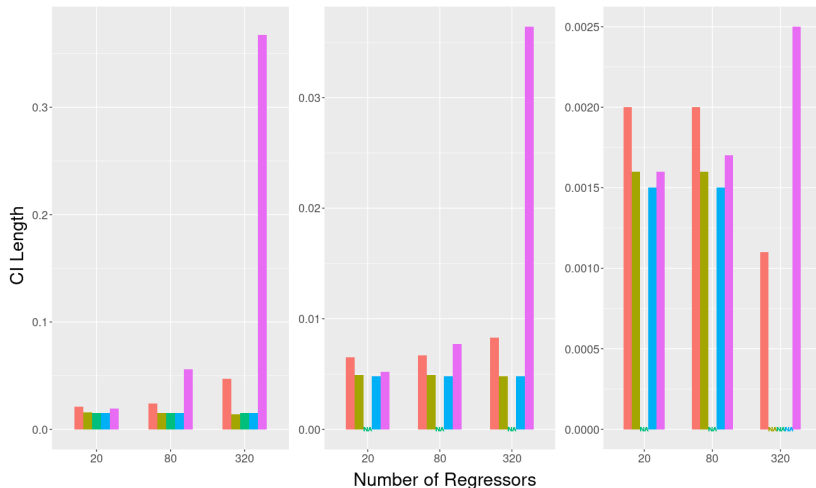


Table: Performance of S-subGD: $n = 10^7$

d	Time (sec.)	Coverage Rate	CI Length
10	5.87	0.965	0.0020
20	11.05	0.955	0.0020
40	21.86	0.954	0.0020
80	43.12	0.952	0.0020
160	81.35	0.953	0.0021
320	166.40	0.963	0.0011
1000	762.16	0.925	0.0461

Application to College Wage Premium

Gender Gap in College Wage Premium

- a stylized fact that the higher college wage premium for women as the major cause for attracting more women to attend and graduate from colleges than men (e.g., Goldin et al. (2006); Chiappori et al. (2009)).
- Current Population Survey (CPS) data : top coded wage issue
⇒ Hubbard's (2011) quantile regression
- Our Goals:
 - ① identify (if any) the heterogeneous effects across quantiles
 - ② properly control other observable characteristics, such as work experience.
 - ③ understand the trends in the college wage premium respectively for female and male
 - ④ understand the difference in gender trends in the college wage premium.

Data I

- We use the samples over six different years (1980, 1990, 2000-2015) from IPUMS USA at <https://usa.ipums.org/usa/>.
- In the years from 1980 to 2000, we use the 5% State sample which is a 1-in-20 national random sample of the population. In the remaining years, we use the American Community Survey (ACS) each year.
- The sampling ratio varies from 1-in-261 to 1-in-232 in 2001-2004, but it is set to a 1-in-100 national random sample of the population after 2005.
- To balance the sample size, we bunch the sample every 5 year after 2001.

Data II

- We restrict our sample to *White*, $18 \leq \text{Age} \leq 65$, and $\text{Wage} \geq \$62$, which is a half of minimum wage earnings in 1980 ($\$3.10 \times 40\text{hours} \times 1/2$).
- *Wage* denotes the implied weekly wage that is computed by dividing yearly earnings by weeks worked last year.
- We only consider full-time workers who worked more than 30 hours per week.
- Then, we compute the *real* wage using the personal consumption expenditures price index (PCEPI) normalized in 1980.
- The data cleaning leaves us with 3.6-4.7 million observations besides 2001-2005, where we have around 2.5 million observations.
- *Educ* denotes an education dummy for some college or above.

Data III

- For control, we use 12 age group dummies with a four-year interval, 51 states dummies (including D.C.), and their interactions. The model contains 1226 covariates in total. We also add 4 additional year dummies for the 5-year combined samples after 2001.

Table: Summary Statistics

Year	Sample Size	$\mathbb{E}(F)$	$\mathbb{E}(Edu)$	$\mathbb{E}(Edu M)$	$\mathbb{E}(Edu F)$
1980	3,659,684	0.390	0.433	0.444	0.416
1990	4,192,119	0.425	0.543	0.537	0.550
2000	4,479,724	0.439	0.600	0.578	0.629
2001-2005	2,493,787	0.447	0.642	0.619	0.670
2006-2010	4,708,119	0.447	0.663	0.631	0.701
2011-2015	4,542,874	0.447	0.686	0.646	0.735

Figure: College Wage Premium: Combining 5-Year Data

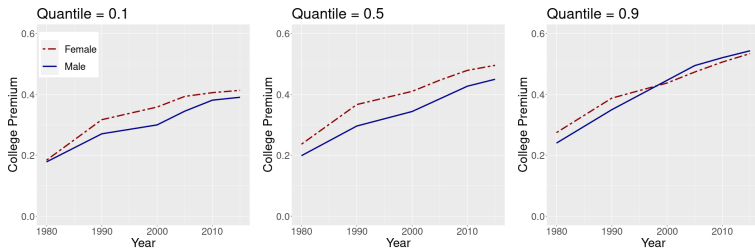


Table: College Wage Premium: $\tau = 0.5$

Year	Female	Male	Difference	Time (min.)
<u>$\tau = 0.5$</u>				
1980	0.2365 [0.2294,0.2435]	0.1988 [0.1945,0.2030]	0.0377 [0.0291,0.0463]	29.4
1990	0.3667 [0.3603,0.3732]	0.2962 [0.2942,0.2982]	0.0705 [0.0634,0.0777]	34.2
2000	0.4101 [0.4056,0.4146]	0.3439 [0.3372,0.3506]	0.0662 [0.0552,0.0772]	36.7
2001-2005	0.4468 [0.4369,0.4567]	0.3854 [0.3765,0.3944]	0.0613 [0.0554,0.0673]	20.2
2006-2010	0.4791 [0.4748,0.4834]	0.4271 [0.4174,0.4368]	0.0520 [0.0454,0.0585]	47.7
2011-2015	0.4957 [0.4887,0.5027]	0.4498 [0.4455,0.4542]	0.0458 [0.0348,0.0568]	46.0

Conclusion

- We provide a new scalable on-line inference method for Quantile Regression with “ultra-large” sample sizes.
- Fast + small memory cost
- Potential Extensions:
 - cluster robust inference
 - high-dimensional settings and penalized estimation like Lasso.
 - more sophisticated random scaling

thank
you