

The copula-graphic estimator in censored nonparametric location-scale regression models

Aleksandar Sujica, Ingrid Van Keilegom

Université catholique de Louvain

YRD 2013, February 1, LLN

Outline

- 1 Introduction (Model)
- 2 Estimator construction
- 3 Simulations
- 4 Further research

Outline

- 1 Introduction (Model)
- 2 Estimator construction
- 3 Simulations
- 4 Further research

- Survival analysis
- Informative censoring
- Copulas
- Location-scale regression

Survival analysis

Example

Y = time to first relapse of a cancer (survival time)

X = age of patient

How to estimate $F(y|x)$?

Survival analysis

Example

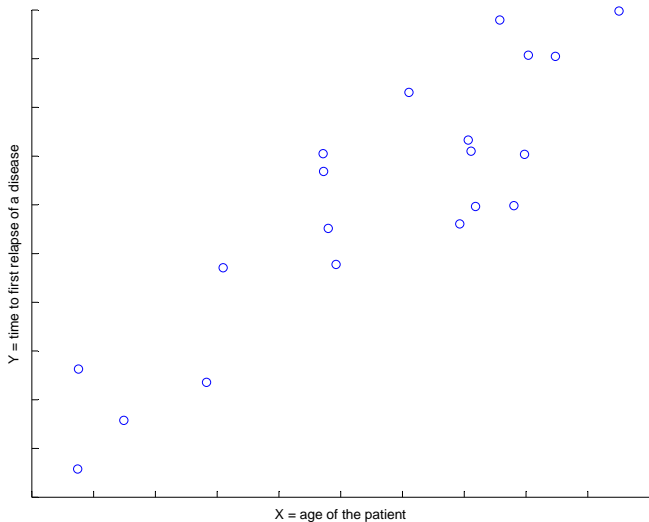
Y = time to first relapse of a cancer (survival time)

X = age of patient

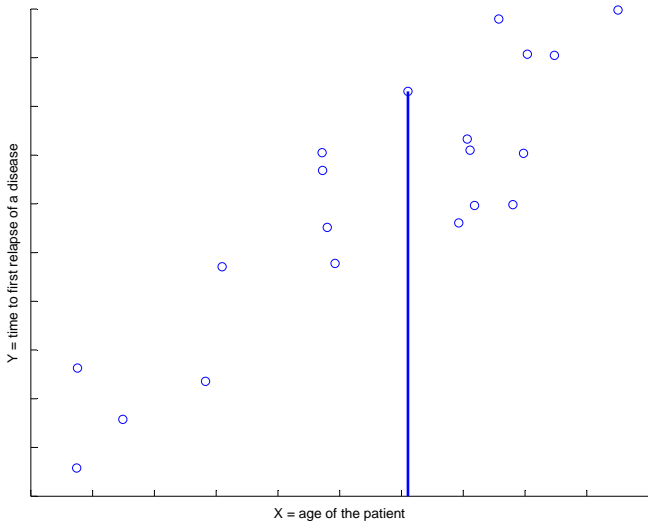
How to estimate $F(y|x)$?

- Survival analysis
- Informative censoring
- Copulas
- Location-scale regression

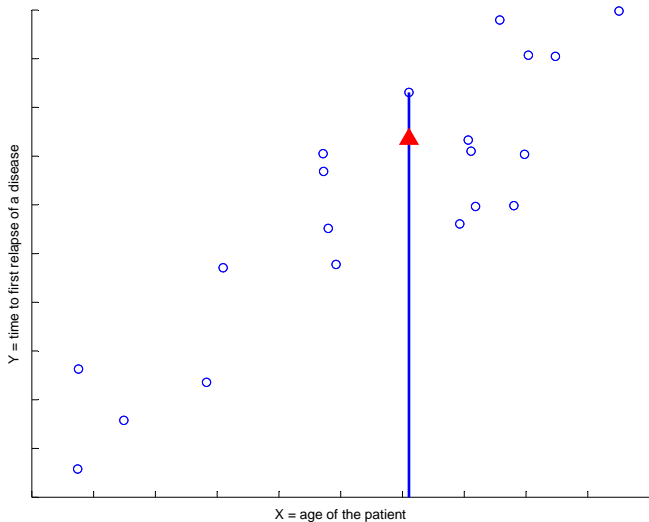
Censoring



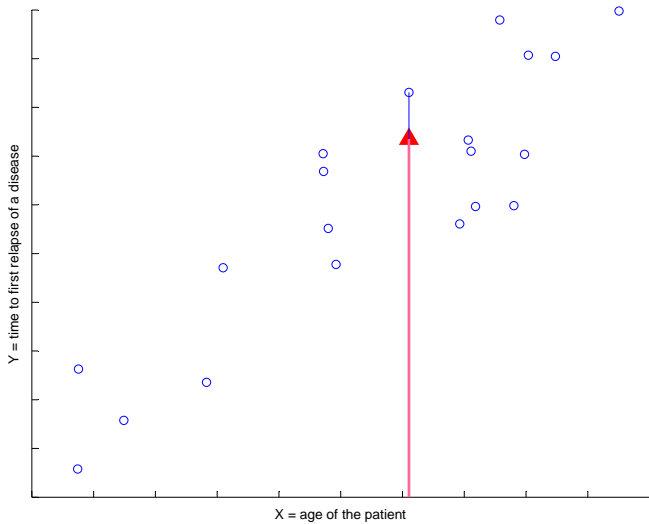
Censoring



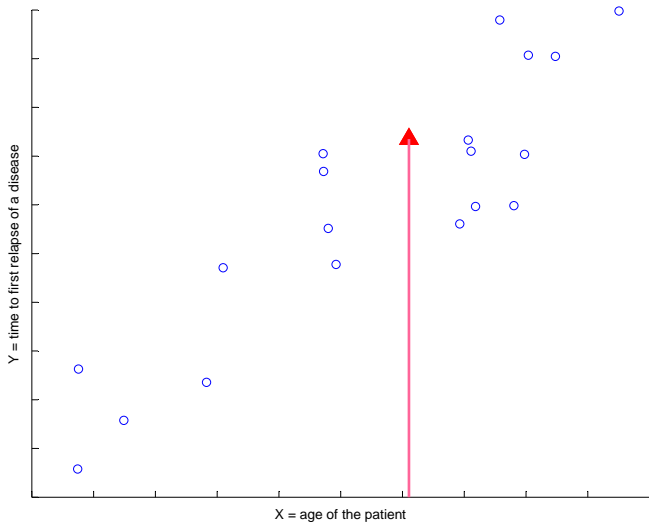
Censoring



Censoring



Censoring



Censoring

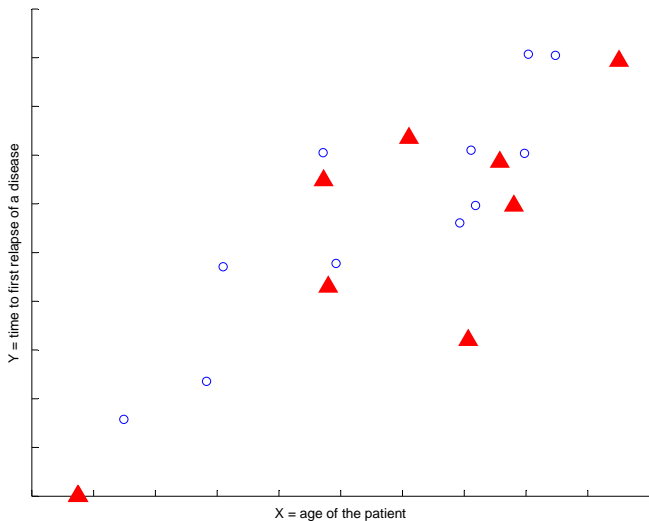
What is Censoring ?

⇒ Instead of observing Y , we observe (T, Δ) , where

- $T = \min(Y, C)$
- $\Delta = I(Y \leq C)$,

and C is the censoring time.

Informative censoring



Copula

What is a copula ?

A function containing ALL INFORMATION about DEPENDENCE between two random variables

Sklar's theorem

If F and G are continuous, there exists a unique copula \mathcal{C} such that

$$P(Y > y, C > c) = \mathcal{C}(\bar{F}(y), \bar{G}(c)),$$

where $\bar{F}(y) = 1 - F(y) = P(Y > y)$

$$\bar{G}(c) = 1 - G(c) = P(C > c)$$

Copula

What is a copula ?

A function containing ALL INFORMATION about DEPENDENCE between two random variables

Sklar's theorem

If F and G are continuous, there exists a unique copula \mathcal{C} such that

$$P(Y > y, C > c) = \mathcal{C}(\bar{F}(y), \bar{G}(c)),$$

where $\bar{F}(y) = 1 - F(y) = P(Y > y)$

$$\bar{G}(c) = 1 - G(c) = P(C > c)$$

Copula

We will use **Archimedean copula** i.e.

$$\mathcal{C}(u, v) = \varphi^{-1}(\varphi(u) + \varphi(v)), \quad 0 \leq u, v \leq 1,$$

⇒ what gives

$$P(Y > y, C > c) = \varphi^{-1}(\varphi(\bar{F}(y)) + \varphi(\bar{G}(c))).$$

For each patient we also observe a **covariate** X

⇒ We will have for every x a generator φ_x such that

$$P(Y > y, C > c | X = x) = \varphi_x^{-1}(\varphi_x(\bar{F}(y|x)) + \varphi_x(\bar{G}(c|x))),$$

where $\bar{F}(y|x) = P(Y > y | X = x)$

$$\bar{G}(c|x) = P(C > c | X = x)$$

Copula

We will use **Archimedean copula** i.e.

$$\mathcal{C}(u, v) = \varphi^{-1}(\varphi(u) + \varphi(v)), \quad 0 \leq u, v \leq 1,$$

⇒ what gives

$$P(Y > y, C > c) = \varphi^{-1}(\varphi(\bar{F}(y)) + \varphi(\bar{G}(c))).$$

For each patient we also observe a **covariate** X

⇒ We will have for every x a generator φ_x such that

$$P(Y > y, C > c | X = x) = \varphi_x^{-1}(\varphi_x(\bar{F}(y|x)) + \varphi_x(\bar{G}(c|x))),$$

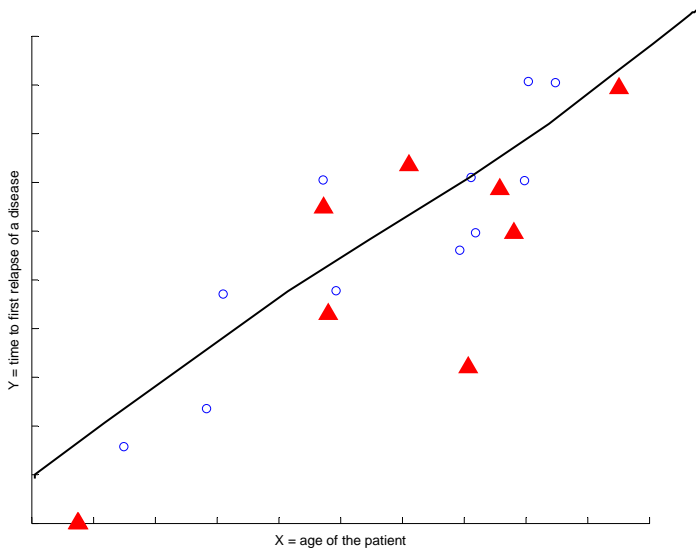
where $\bar{F}(y|x) = P(Y > y | X = x)$

$$\bar{G}(c|x) = P(C > c | X = x)$$

- Survival analysis
- Informative censoring
- Copulas
- Location-scale regression

Location-scale regression

$$Y = m(X) + \sigma(X)\varepsilon, \quad X \perp\!\!\!\perp \varepsilon$$



- Survival analysis
- Informative censoring
- Copulas
- Location-scale regression

Our model

- $(T_i = \min(Y_i, C_i), \Delta_i, X_i) \quad i = 1, \dots, n$
- $P(Y < y, C < c | X = x) = \Phi_x^{-1}(\Phi_x(F(y|x)) + \Phi_x(G(c|x)))$
- $Y = m(X) + \sigma(X)\varepsilon, \quad X \perp\!\!\!\perp \varepsilon$

History and Contribution

Example

Y = time to first relapse of a cancer (survival time)

C = time until leaving the study

X = age of the patient

How to estimate $F(y|x)$?

⇒ Beran (1981)

Y, C given X are independent

⇒ Van Keilegom & Akritas (1999)

Y, C given X are independent

$$Y = m(X) + \sigma(X)\varepsilon, \quad X \perp\!\!\!\perp \varepsilon$$

⇒ Braekers & Veraverbeke (2005)

Y, C given X are copula dependent and copula is known

History and Contribution

Example

Y = time to first relapse of a cancer (survival time)

C = time until leaving the study

X = age of the patient

How to estimate $F(y|x)$?

⇒ Beran (1981)

Y, C given X are independent

⇒ Van Keilegom & Akritas (1999)

Y, C given X are independent

$$Y = m(X) + \sigma(X)\varepsilon, \quad X \perp\!\!\!\perp \varepsilon$$

⇒ Braekers & Veraverbeke (2005)

Y, C given X are copula dependent and copula is known

History and Contribution

Example

Y = time to first relapse of a cancer (survival time)

C = time until leaving the study

X = age of the patient

How to estimate $F(y|x)$?

⇒ Beran (1981)

Y, C given X are independent

⇒ Van Keilegom & Akritas (1999)

Y, C given X are independent

$$Y = m(X) + \sigma(X)\varepsilon, \quad X \perp\!\!\!\perp \varepsilon$$

⇒ Braekers & Veraverbeke (2005)

Y, C given X are copula dependent and copula is known

History and Contribution

Example

Y = time to first relapse of a cancer (survival time)

C = time until leaving the study

X = age of the patient

How to estimate $F(y|x)$?

⇒ Beran (1981)

Y, C given X are independent

⇒ Van Keilegom & Akritas (1999)

Y, C given X are independent

$$Y = m(X) + \sigma(X)\varepsilon, \quad X \perp\!\!\!\perp \varepsilon$$

⇒ Braekers & Veraverbeke (2005)

Y, C given X are copula dependent and copula is known

	Fully nonparametric model	Location-scale regression
Independent censoring	Beran (1981)	VK & Akritas (1999)
Informative censoring	Braekers & Veraverbeke (2005)	Our model

Outline

- 1 Introduction (Model)
- 2 Estimator construction**
- 3 Simulations
- 4 Further research

Estimating $\bar{F}(y|x) = P(Y > y|X = x)$

For simplicity reasons assume that $\sigma(x) = 1$

$$Y = m(X) + \varepsilon$$

We can show that

$$\bar{F}(y|x) = \bar{F}_e(y - m(x)),$$

where $\bar{F}_e(t) = P(\varepsilon > t)$.

We will estimate $\bar{F}(y|x)$ by estimating $\bar{F}_e(\cdot)$ and $m(\cdot)$

$$\hat{\bar{F}}(y|x) = \hat{\bar{F}}_e(y - \hat{m}(x))$$

Estimating $\bar{F}(y|x) = P(Y > y|X = x)$

For simplicity reasons assume that $\sigma(x) = 1$

$$Y = m(X) + \varepsilon$$

We can show that

$$\bar{F}(y|x) = \bar{F}_e(y - m(x)),$$

where $\bar{F}_e(t) = P(\varepsilon > t)$.

We will estimate $\bar{F}(y|x)$ by estimating $\bar{F}_e(\cdot)$ and $m(\cdot)$

$$\hat{\bar{F}}(y|x) = \hat{\bar{F}}_e(y - \hat{m}(x))$$

Estimating $\bar{F}(y|x) = P(Y > y|X = x)$

For simplicity reasons assume that $\sigma(x) = 1$

$$Y = m(X) + \varepsilon$$

We can show that

$$\bar{F}(y|x) = \bar{F}_e(y - m(x)),$$

where $\bar{F}_e(t) = P(\varepsilon > t)$.

We will estimate $\bar{F}(y|x)$ by estimating $\bar{F}_e(\cdot)$ and $m(\cdot)$

$$\hat{\bar{F}}(y|x) = \hat{\bar{F}}_e(y - \hat{m}(x))$$

Estimating $\bar{F}(y|x) = P(Y > y|X = x)$

For simplicity reasons assume that $\sigma(x) = 1$

$$Y = m(X) + \varepsilon$$

We can show that

$$\bar{F}(y|x) = \bar{F}_e(y - m(x)),$$

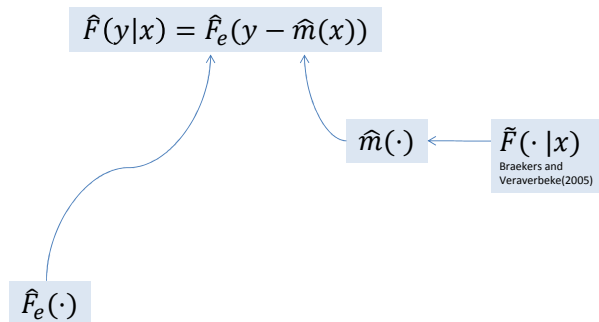
where $\bar{F}_e(t) = P(\varepsilon > t)$.

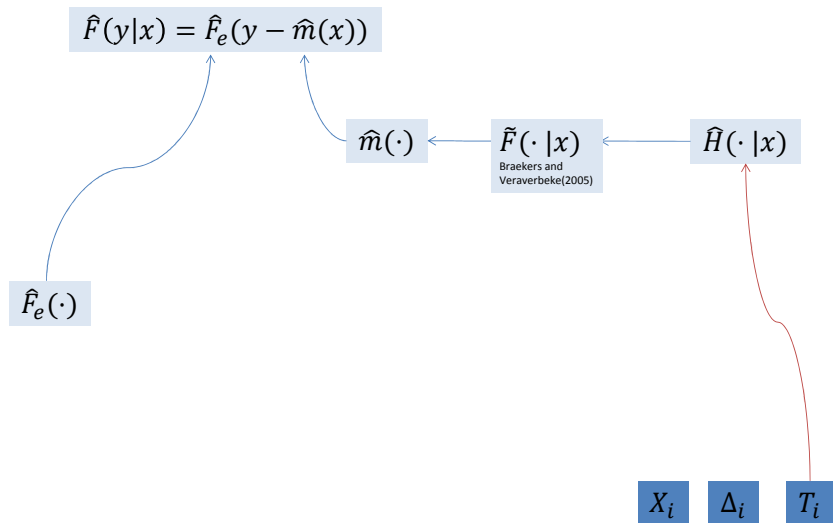
We will estimate $\bar{F}(y|x)$ by estimating $\bar{F}_e(\cdot)$ and $m(\cdot)$

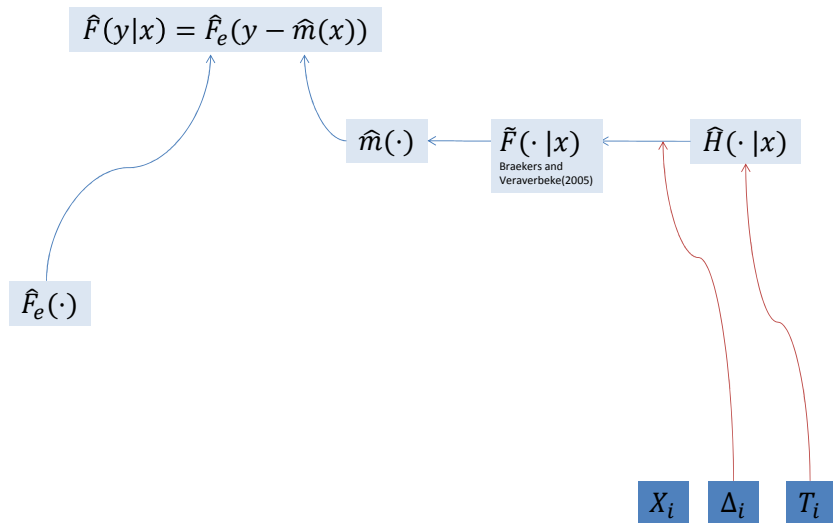
$$\hat{\bar{F}}(y|x) = \hat{\bar{F}}_e(y - \hat{m}(x))$$

$$\hat{F}(y|x) = \hat{F}_e(y - \hat{m}(x))$$

 $\hat{m}(\cdot)$ $\hat{F}_e(\cdot)$ X_i Δ_i T_i

 X_i Δ_i T_i





$$\hat{F}(y|x) = \hat{F}_e(y - \hat{m}(x))$$

$$\hat{m}(\cdot)$$

$$\tilde{F}(\cdot | x)$$

Braekers and
Veraverbeke(2005)

$$\hat{H}(\cdot | x)$$

$$\hat{F}_e(\cdot) = \hat{\phi}_{(t)} \left\{ - \int_B \int_{-\infty}^{(\cdot)} \phi'_x(\hat{H}_e(s|x)) d\hat{H}_e^u(s|x) d\hat{F}_X(x) \right\}$$

 X_i
 Δ_i
 T_i

$$\hat{F}(y|x) = \hat{F}_e(y - \hat{m}(x))$$

$$\hat{m}(\cdot)$$

$$\tilde{F}(\cdot |x)$$

Braekers and
Veraverbeke(2005)

$$\hat{H}(\cdot |x)$$

$$\hat{F}_e(\cdot) = \hat{\phi}(t) \left\{ - \int_B \int_{-\infty}^{(\cdot)} \phi'_x(\hat{H}_e(s|x)) d\hat{H}_e^u(s|x) d\hat{F}_X(x) \right\}$$

$$\hat{H}_e(\cdot |x) = \sum_{i=1}^n W_{ni}(x, h_n) I\left(\frac{T_i - \hat{m}(X_i)}{\sigma(X_i)} \leq \cdot\right)$$

 X_i
 Δ_i
 T_i

$$\hat{F}(y|x) = \hat{F}_e(y - \hat{m}(x))$$

$$\hat{m}(\cdot)$$

$$\tilde{F}(\cdot | x)$$

Braekers and
Veraverbeke(2005)

$$\hat{H}(\cdot | x)$$

$$\hat{F}_e(\cdot) = \hat{\phi}(t) \left\{ - \int_B \int_{-\infty}^{(\cdot)} \phi'_x(\hat{H}_e(s|x)) d\hat{H}_e^u(s|x) d\hat{F}_X(x) \right\}$$

$$\hat{H}_e(\cdot | x) = \sum_{i=1}^n W_{ni}(x, h_n) I\left(\frac{T_i - \hat{m}(X_i)}{\sigma(X_i)} \leq \cdot\right)$$

 X_i
 Δ_i
 T_i

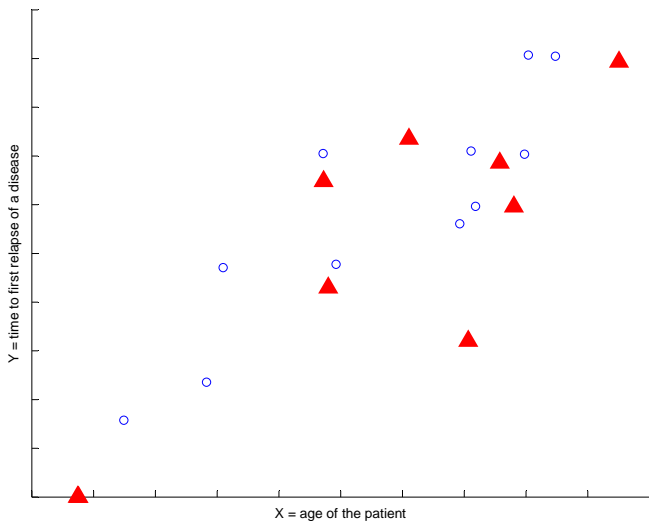
Outline

- 1 Introduction (Model)
- 2 Estimator construction
- 3 Simulations**
- 4 Further research

	Fully nonparametric model	Location-scale regression
Independent censoring	Beran (1981)	VK & Akritas (1999)
Informative censoring	Braekers & Veraverbeke (2005)	New estimator

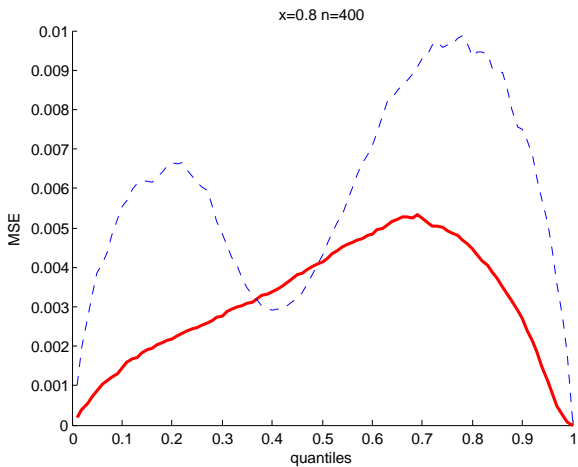
Simulation framework

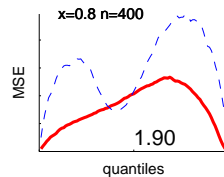
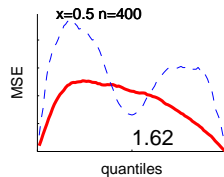
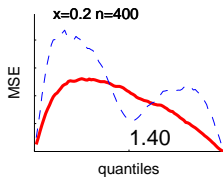
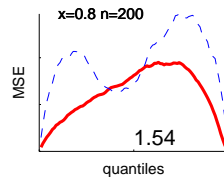
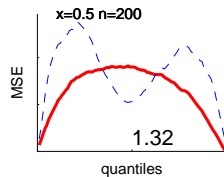
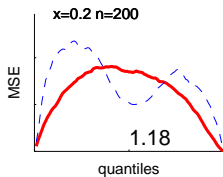
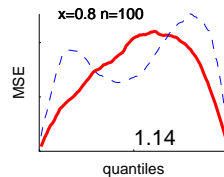
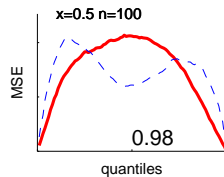
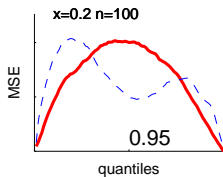
$$Y = 0.1X + 0.015\varepsilon$$



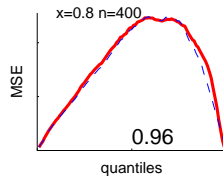
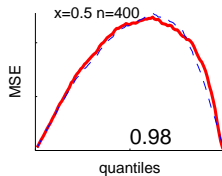
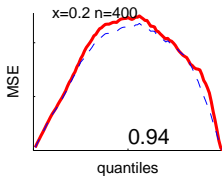
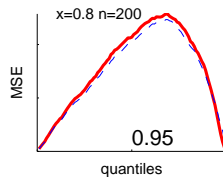
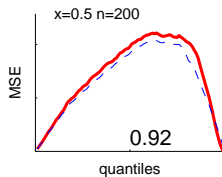
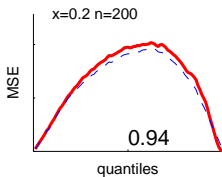
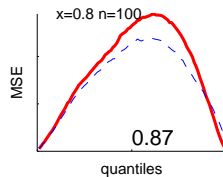
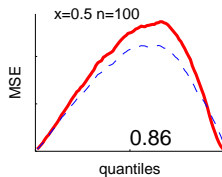
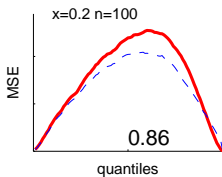
New estimator

Braekers & Veraverbeke (2005)





	Fully nonparametric model	Location-scale regression
Independent censoring	Beran (1981)	VK & Akritas (1999)
Informative censoring	Braekers & Veraverbeke (2005)	New estimator



Outline

- 1 Introduction (Model)
- 2 Estimator construction
- 3 Simulations
- 4 Further research**

We have studied the estimation of $F(y|x) = P(Y \leq y|X = x)$ when

- ◇ $Y = m(X) + \sigma(X)\varepsilon$, with ε independent of X
- ◇ Y is subject to random right censoring
- ◇ Y and C are copula dependent for given X

Future research :

- ◇ Tests for the comparison of regression curves
- ◇ Tests for the error distribution $F_\varepsilon(\cdot)$
- ◇ Estimation of copula

Thank you for your time :)

Black guy

