

# Semi-supervised Algorithms for Risk Assessment with Noisy EHR Data

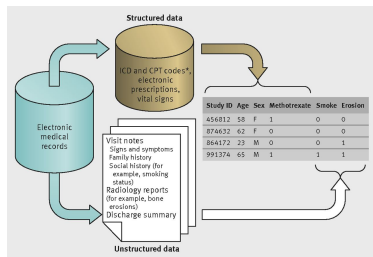
Chuan Hong

(Joint work with Tianxi Cai from Harvard University)

Department of Biostatistics and Bioinformatics  
Duke University

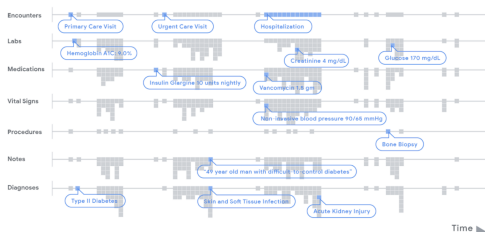
September 8, 2022

# EHR Structure



(Figure from Liao et al. 2015)

Patient Timeline



(Figure from ai.googleblog.com)

- ▶ **Structured data:** ICD billing codes; lab results etc
- ▶ **Unstructured text data:** extracted via natural language processing (NLP)
- ▶ **Detailed longitudinal patient level data**

# Challenges in EHR-linked Survival Analysis

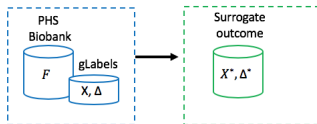
Disease status and event time information is not readily available.

- ▶ Annotating **event time** requires labor extensive chart review
- ▶ Time to first ICD codes **inaccurate**
- ▶ **Surrogate event time**: derived from label+codes+NLP
- ▶  $\rightsquigarrow$  power loss, biased estimates

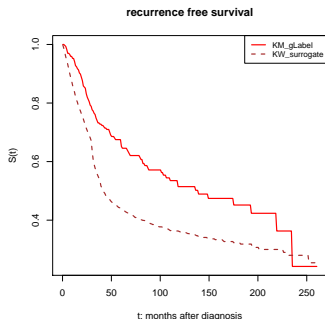
# Motivating Example: PHS Lung Cancer Study

- ▶ PHS: Partners Healthcare contains both a wealth of clinical and also biological measurements.

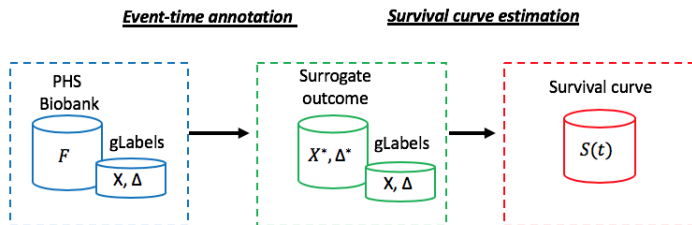
## Event-time annotation



- ▶ Aim to estimate recurrence free survival for lung cancer patients.
- ▶ Lung cancer data mart
- ▶ 70K patients from PHS biobank EMR
- ▶ 40K patients identified as lung cancer
- ▶ 5K early stage patients
- ▶ 300 gLabels for event time manually annotated by domain experts
- ▶ Surrogate outcomes derived by using Uno et al. 2018



# Motivating Example: PHS Lung Cancer Study



Solution: Develop a **calibrated survival curve** that combines imperfect sources of information on event time in EHR, together with the exact event time.

# Sources of information on event times in EHRs

- ▶  $T_i$ : failure time, the time that the patient developed the event
- ▶  $C_i$ : censoring time

Labeled data:  $i = 1, \dots, n$

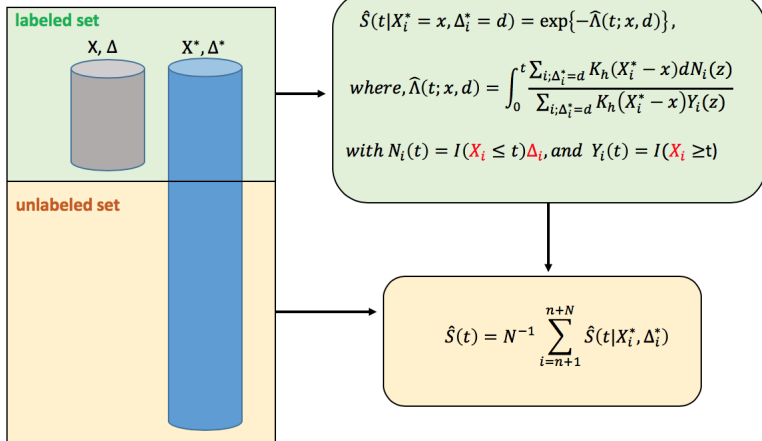
- ▶  $X_i = \min(T_i, C_i)$ : observed event time:
- ▶  $\Delta_i = \mathbb{1}(T_i \leq C_i)$ : the censoring indicator (whether the patient developed the event prior to the last visit)

Unlabeled data:  $i = 1, \dots, n + N$

- ▶  $X^*$ : imperfect estimates of event times
  - time to the first ICD code related to the event
  - time to the first NLP mention of the event
  - algorithm annotated event time
- ▶  $\Delta_i^* = \mathbb{1}(T_i^* \leq C_i)$

Goal: estimate the survival function  $S(t)$ .

# Conditional Nelson-Aalen Estimator



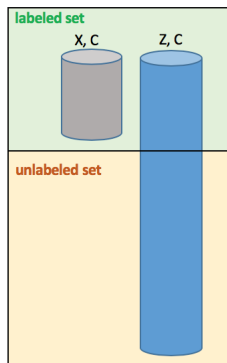
- ▶ Parast et al. 2014: baseline covariates
- ▶ In our case,  $T$  and  $C$  may no longer be independent conditional on  $X^*$  and  $\Delta^*$  if there is no restriction on  $X^*$  and  $\Delta^*$ .
- ▶ As a result, the Nelson-Aalen type estimator defined above may not be consistent.

# Semiparametric Estimator

Let  $\mathbf{Z}$  denote a variable related to both  $T$  and  $C$ .

As  $T \perp\!\!\!\perp C$ , we have

$$\begin{aligned}\pi_t &= E\{\mathbb{1}(T \leq t)\} = E\{\mathbb{1}(T \leq t) \mid C > t\} = E\{\mathbb{1}(T \leq t)\Delta \mid C > t\} \\ &= E[E\{N(t) \mid \mathbf{Z}, C > t\} \mid C > t].\end{aligned}$$



**Step 1: fit a time-varying coefficient working model**

Let  $N(t) = I(T \leq t \wedge C)$ , and  $N_i(t) = I(T_i \leq t \wedge C_i)$

$$E\{N(t) \mid \mathbf{Z}, C > t\} = g(\alpha_t + \beta_t' \mathbf{Z}_i) = \frac{\exp(\alpha_t + \beta_t' \mathbf{Z}_i)}{1 + \exp(\alpha_t + \beta_t' \mathbf{Z}_i)}$$

**Step 2: derive an estimate of the survival curve**

$$\hat{S}(t) = 1 - \frac{\sum_{i=n+1}^{n+N} I(C_i > t) g(\hat{\alpha}_t + \hat{\beta}_t' \mathbf{Z}_i)}{\sum_{i=n+1}^{n+N} I(C_i > t)}$$



# Combined Estimator

To improve the efficiency, we combine proposed semiparametric estimator with the KM.

Let  $\hat{\boldsymbol{\mu}} = (\hat{\mathbf{S}}_{\text{Semi}}, \hat{\mathbf{S}}_{\text{KM}})^{\top}$ , and  $\boldsymbol{\Sigma}$  denote their covariance matrix, then the combined estimator is constructed as

$$(\mathbf{1}^{\top} \boldsymbol{\Sigma}^{-1} \mathbf{1})^{-1} \mathbf{1}^{\top} \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\mu}}$$

- ▶ unbiased
- ▶ smallest possible variance among all linear combinations

# Asymptotic Properties of Proposed Estimator

Let  $\mathbf{W}_i = (1, \mathbf{z}_i^T)^T$ .

Step 1.

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_t) \rightarrow \text{Normal}\{\mathbf{0}, \mathbf{A}^{-1}\mathbf{B}(\mathbf{A}^{-1})^T\},$$

where

$$\mathbf{A} = -E\{\mathbb{1}(C > t)g'(\boldsymbol{\theta}_t^T \mathbf{W})\mathbf{W}\mathbf{W}^T\};$$

$$\mathbf{B} = \text{cov}[\mathbb{1}(C > t)\{N(t) - g(\boldsymbol{\theta}_t^T \mathbf{W})\}\mathbf{W}];$$

$$g'(x) = \exp(x)/\{1 + \exp(x)\}^2.$$

Step 2.

$$\sqrt{n}(\hat{\pi}_t - \pi_t) = G(t)^{-1}E\{\mathbb{1}(C > t)g'(\boldsymbol{\theta}_t^T \mathbf{W})\mathbf{W}\}^T \sqrt{n}(\hat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_t) + o_p(1),$$

where  $G(t) = P(C > t)$ .

# Asymptotic Properties of Kaplan Meier Estimator

We have

$$\sqrt{n}(\hat{\pi}_t^{\text{KM}} - \pi_t) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (1 - \pi_t) \int_0^t \frac{dM_i(u)}{P(X \geq u)} + o_p(1).$$

In practice, we have to replace all the unknown quantities in the above influence function with their estimations. This leads to

$$(1 - \hat{\pi}_t^{\text{KM}}) \int_0^t \frac{d\hat{M}_i(u)}{\bar{Y}(u)} = (1 - \hat{\pi}_t^{\text{KM}}) \left\{ \frac{\Delta_i \mathbb{1}(X_i \leq t)}{\bar{Y}(X_i)} - \sum_j \frac{\Delta_j \mathbb{1}(X_j \leq t \wedge X_j)}{n\bar{Y}^2(X_j)} \right\},$$

where  $\bar{Y}(u) = n^{-1} \sum Y_i(u)$  and  $Y_i(u) = \mathbb{1}(X_i \geq u)$ .

# Combined Estimator

To improve the efficiency, we combine proposed semiparametric estimator with the KM.

Let  $\hat{\boldsymbol{\mu}} = (\hat{\mathbf{S}}_{\text{Semi}}, \hat{\mathbf{S}}_{\text{KM}})^{\top}$ , and  $\boldsymbol{\Sigma}$  denote their covariance matrix, then the combined estimator is constructed as

$$(\mathbf{1}^{\top} \boldsymbol{\Sigma}^{-1} \mathbf{1})^{-1} \mathbf{1}^{\top} \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\mu}}$$

where

$$\Sigma_{1,1} = \mathbf{C}^{\top} E[\mathbb{1}(C_j > t) \{N_j(t) - g(\boldsymbol{\theta}_t^{\top} \mathbf{W}_j)\}^2 \mathbf{W}_j \mathbf{W}_j^{\top}] \mathbf{C}.$$

$$\Sigma_{2,2} = E[\mathbb{1}(T_i \wedge t \leq C_i) \{\pi_t^{\text{KM}} - \mathbb{1}(X_i < t)\}^2 / G^2(X_i \wedge t)].$$

$$\Sigma_{1,2} = -\mathbf{C}^{\top} E[\mathbb{1}(C_i > t) \{N_i(t) - g(\boldsymbol{\theta}_t^{\top} \mathbf{W}_i)\} \mathbf{W}_i w_i(t) \{\pi_t^{\text{KM}} - \mathbb{1}(X_i < t)\}]$$

- ▶ unbiased
- ▶ smallest possible variance among all linear combinations

# Simulation Setup

- ▶ 400 datasets, each has 200 labeled subjects and 2000 unlabeled subjects
- ▶ For each dataset, generate
  - ▶  $T_i \sim \text{Exponential}(1)$
  - ▶  $C_i \sim \text{Exponential}(1)$  or  $C_i \sim \text{Uniform}(3)$
  - ▶  $Z_i = h(\lambda, T_i, C_i) + e_i$ , with  $e_i \sim \text{Normal}(0, \sigma)$
  - ▶  $\lambda$  large vs. small

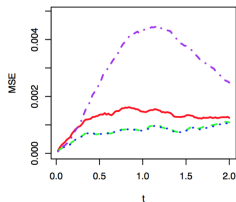
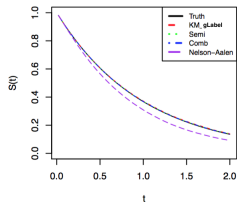
# Settings

1.  $T_i \sim \text{Exponential}(1)$ ,  $C_i \sim \text{Exponential}(1)$ ,  
 $Z_i = \log(T_i) + \lambda \log(C_i) + e_i$ , with  $e_i \sim \text{Normal}(0, 1)$ ;  
 $\lambda = 1$  and  $\lambda = 0.1$ .
2.  $T_i \sim \text{Exponential}(1)$ ,  $C_i \sim \text{Exponential}(1)$ ,  
 $Z_i = \lambda \log(T_i) + (1 - \lambda) \log(C_i) + e_i$ , with  $e_i \sim \text{Normal}(0, 1)$ ;  
 $\lambda = 1/2$  versus  $\lambda = 0.9$ .
3.  $T_i \sim \text{Exponential}(1)$ ,  $C_i \sim \text{Exponential}(1)$ ,  
 $Z_i = \lambda \log(T_i) + (1 - \lambda) \log(C_i) + e_i$ , with  $e_i \sim \text{Normal}(0, 0.25)$ ;  
 $\lambda = 1/2$  versus  $\lambda = 0.9$ .
4.  $T_i \sim \text{Exponential}(1)$ ,  $C_i \sim \text{Uniform}(0, 3)$ ,  
 $Z_i = \lambda \log(T_i) + (1 - \lambda) \log(C_i) + e_i$ , with  $e_i \sim \text{Normal}(0, 0.25)$ ;  
 $\lambda = 1/2$  versus  $\lambda = 0.9$ .
5.  $T_i \sim \text{Exponential}(1)$ ,  $C_i \sim \text{Uniform}(0, 3)$ ,  
 $Z_i = \log\{\min(T_i^*, C_i)\}$ , with  $T_i^* = T_i + e_i$  and  $e_i \sim \text{Exponential}(\lambda)$ ;  
 $\lambda = 2$  versus  $\lambda = 5$ .

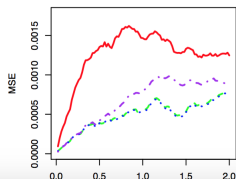
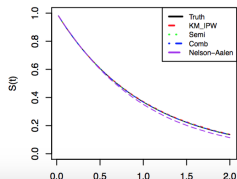
# Simulation Results

$T_i \sim \text{Exponential}(1)$ ,  $C_i \sim \text{Uniform}(0, 3)$ ,  
 $Z_i = \lambda \log(T_i) + (1 - \lambda) \log(C_i) + e_i$ , with  $e_i \sim \text{Normal}(0, 0.25)$ ;  
 $\lambda = 1/2$  versus  $\lambda = 0$

more error



less error



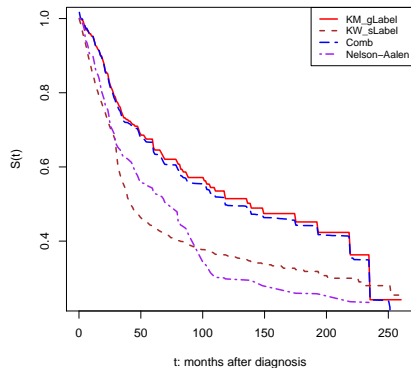
# Real data example

- ▶ Goal: estimate the recurrence-free survival curve of lung cancer patients
- ▶ Dataset
  - ▶ 37,021 total lung cancer patients
  - ▶ 5K early stage patients
  - ▶ 340 had recurrence status and observed time to recurrence labels from chart review ( $\Delta_i, X_i$ )
  - ▶ Surrogate outcomes ( $\Delta_i^*, Z_i$ ) are available for each patient obtained by Uno's two-step estimator
    - binary recurrence status
    - predict event times using peaks of specific features
  - ▶ Accuracy of  $\Delta_i^*$  is only 0.79;

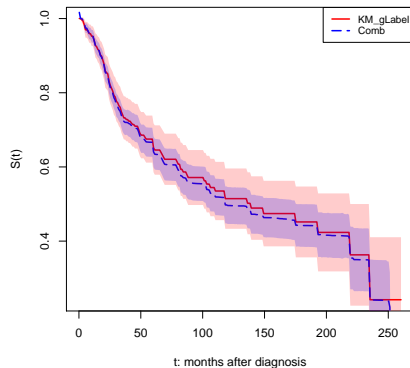


# Results - Survival Curve Comparisons

recurrence free survival



recurrence free survival with CI



## Remarks

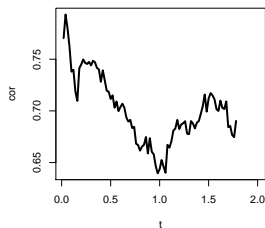
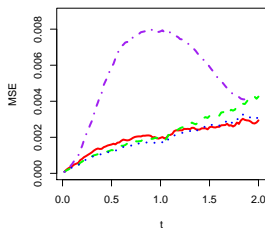
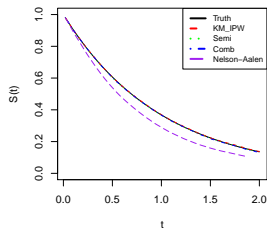
- ▶ a semi-supervised calibrated survival curve
- ▶ fully utilize both labeled and unlabeled data
- ▶ next step: estimate the survival function among different risk groups (e.g., different treatment group), and test for difference.

Thank you!

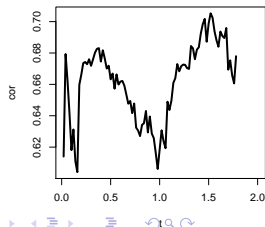
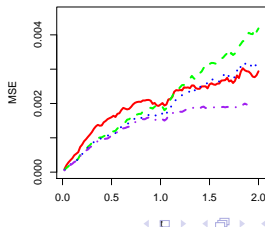
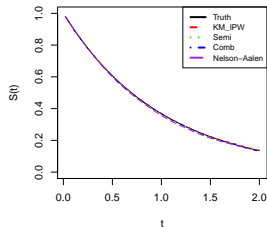
# Setting 1

$T_i \sim \text{Exponential}(1)$ ,  $C_i \sim \text{Exponential}(1)$ ,  
 $Z_i = \log(T_i) + \lambda \log(C_i) + e_i$ , with  $e_i \sim \text{Normal}(0, 1)$ ;  
 $\lambda = 1$  and  $\lambda = 0.1$ .

more error



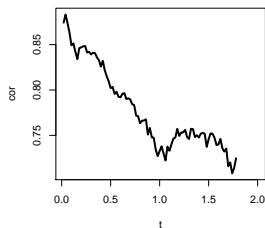
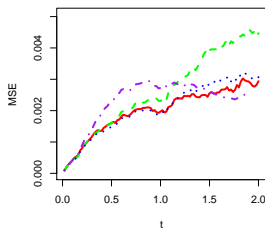
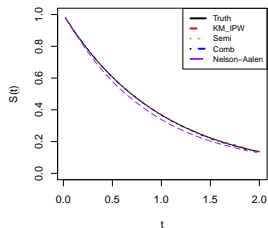
less error



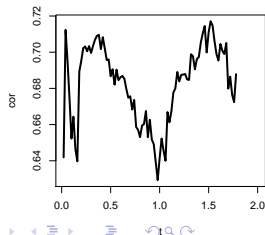
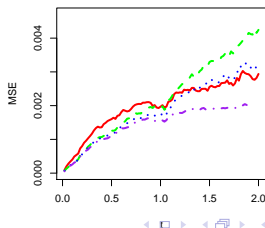
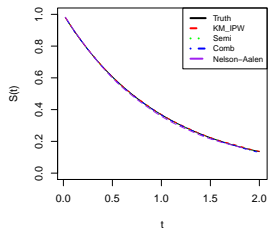
## Setting 2

$T_i \sim \text{Exponential}(1)$ ,  $C_i \sim \text{Exponential}(1)$ ,  
 $Z_i = \lambda \log(T_i) + (1 - \lambda) \log(C_i) + e_i$ , with  $e_i \sim \text{Normal}(0, 1)$ ;  
 $\lambda = 1/2$  versus  $\lambda = 0.9$ .

more error



less error



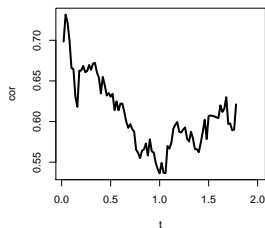
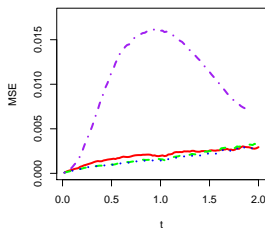
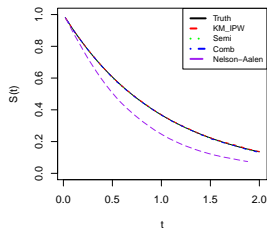
# Setting 3

$T_i \sim \text{Exponential}(1)$ ,  $C_i \sim \text{Exponential}(1)$ ,

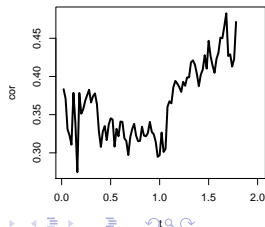
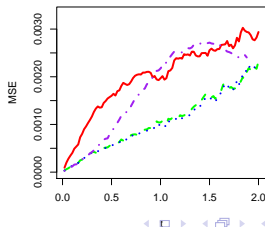
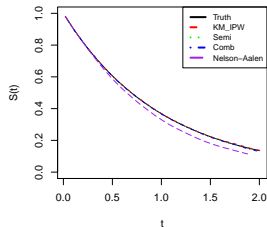
$Z_i = \lambda \log(T_i) + (1 - \lambda) \log(C_i) + e_i$ , with  $e_i \sim \text{Normal}(0, 0.25)$ ;

$\lambda = 1/2$  versus  $\lambda = 0.9$ .

more error



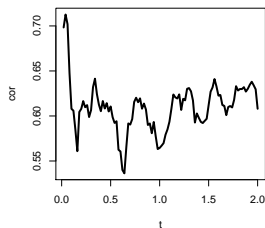
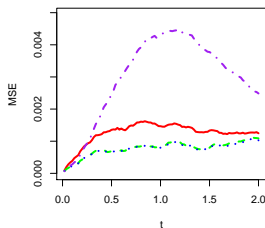
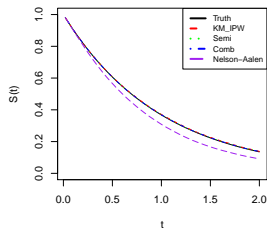
less error



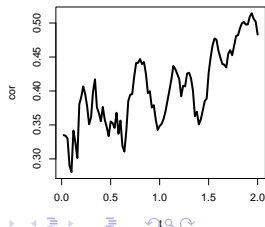
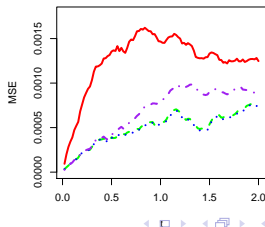
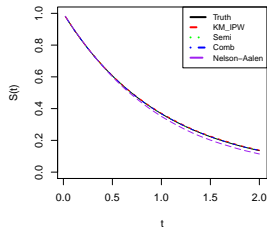
## Setting 4

$T_i \sim \text{Exponential}(1)$ ,  $C_i \sim \text{Uniform}(0, 3)$ ,  
 $Z_i = \lambda \log(T_i) + (1 - \lambda) \log(C_i) + e_i$ , with  $e_i \sim \text{Normal}(0, 0.25)$ ;  
 $\lambda = 1/2$  versus  $\lambda = 0.9$ .

more error



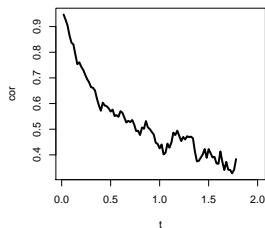
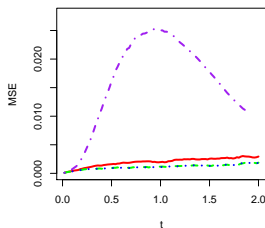
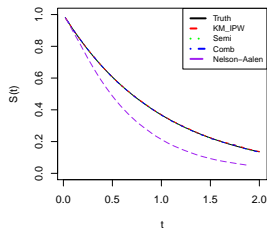
less error



## Setting 5

$T_i \sim \text{Exponential}(1)$ ,  $C_i \sim \text{Uniform}(0, 3)$ ,  
 $Z_i = \log\{\min(T_i^*, C_i)\}$ , with  $T_i^* = T_i + e_i$  and  $e_i \sim \text{Exponential}(\lambda)$ ;  
 $\lambda = 2$  versus  $\lambda = 5$ .

more error



less error

