

Cutpoint Determination Method for the Survival Data with Competing Risks

Jinheum Kim¹, Sook-young Woo², Seonwoo Kim², & JoonSuk Park³

¹Department of Applied Statistics, University of Suwon, Korea

²Biostatistics team, Samsung Biomedical Research Institute, Samsung Advanced Institute of Technology, Korea

³Department of Thoracic Surgery, Samsung Medical Center, Sungkyunkwan University School of Medicine, Korea

August 30, 2012

Outline

- Review on outcome-oriented approaches for survival data
- Propose an approach for identifying a cutpoint & significance test
- Simulations
- Application to the lung cancer data
- Concluding remarks

Why discretizing a continuous covariate?

- Simplest interpretation
 - eg, in Cox model with a binary covariate, $\exp(\beta)$, interpreted as the relative risk
- Interested in identification of a therapeutic threshold for clinical use and treatment decisions

Two approaches

- Based on historical data or or based on a split into groups at a predetermined percentile of the continuous covariate
⇒ “Data-oriented approach”
- Split into groups based on either the largest value of the likelihood or the largest value of two-sample test statistic after a search of possible cutpoints
⇒ “Outcome-oriented approach”
 - Possibly lead to an inflated type-I error due to multiple testing
 - Correction needs to be applied to obtain the correct type-I error

Review: outcome-oriented approach

- Jespersen (1986): base an adjusted test on the maximum value of the score statistic
- Contal & O'Quigley (1999): modify the log-rank test statistic
- Lausen & Schumacher (1992, 1996): propose a standardized test statistic for the two-sample problem with groups defined by a threshold parameter
- Klein & Wu (2004): extend the Contal & O'Quigley's approach to the AFT model and the Cox model with additional covariates

Motivating data

- The 1965 lung cancer patients received a tumor removal surgery at Seoul Samsung Hospital in Korea from September 1994 to December 2005
- Event of interest: relapse, competing risk(s): death, prognostic factor: tumor size at surgery
- Want to split the patients into two groups such as a high risk group and a low risk group to relapse to apply different treatment to each group
⇒ extend the Contal & O'Quigley's approach to the competing risks model through the Gray's statistic (1988)

Notation

- Observed data: $\{(T_i = X_i \wedge Y_i \wedge C_i, \delta_i, Z_i), i = 1, \dots, n\}$
 - X_i (Y_i , or C_i) : event time of interest(time of competing risks or censoring time)
 - δ_i : 1(event occurred), 2(competing risks occurred), 0(censored)
 - Z_i : prognostic factor
- $t_{(1)} < t_{(2)} < \dots < t_{(k)}$: times when an event of interest occurred
- t_{i1}, \dots, t_{im_i} : censored times or failed times due to competing risks in $[t_{(i)}, t_{(i+1)})$, $i = 0, 1, \dots, k$, when $t_{(0)} = 0$ and $t_{(k+1)} = \infty$
- d_{gi} : the number of events of interest in group g at time $t_{(i)}$, $g = 0, 1$; $i = 1, \dots, k$
 - $d_i = d_{0i} + d_{1i}$
- r_{gi} : the number at risk in group g at time $t_{(i)}$
 - $r_i = r_{0i} + r_{1i}$

Gray's statistic

- $F_g(t)$: cumulative incidence function(CIF) for the event of interest in group g at time t
- $S_g(t)$: survival function of being free of any event in group g at time t
- $\hat{F}_g(t) = \sum_{i:t_{(i)} \leq t} \hat{S}_g(t_{(i-1)}) \frac{d_{gi}}{n_{gi}}$: estimated CIF in group g at time t
- $w_{gi} = \frac{1 - \hat{F}_g(t_{(i-1)})}{\hat{S}_g(t_{(i-1)})}$: correction factor in group g at time $t_{(i)}$
- $\tilde{r}_{gi} = \frac{1 - \hat{F}_g(t_{(i-1)})}{\hat{S}_g(t_{(i-1)})} r_{gi} = w_{gi} r_{gi}$: adjusted r_{gi}
 - $\tilde{r}_i = \tilde{r}_{0i} + \tilde{r}_{1i}$
- $\tilde{U}_1 = \sum_{i=1}^k (d_{1i} - d_i \frac{\tilde{r}_{1i}}{\tilde{r}_i})$: Gray's statistic in two-sample case

Linear rank statistic in competing risks model

- $\Phi_1^{(i)}, \dots, \Phi_{d_i}^{(i)}$: 0 if each subject failed at $t_{(i)}$ belongs to group 0, and 1 otherwise
- $\Phi_{i1}, \dots, \Phi_{im_i}$: 0 if each subject censored or failed due to competing risks in $[t_{(i)}, t_{(i+1)})$ belongs to group 0, and 1 otherwise
- For $i = 1, \dots, k$, let

$$a_l^{(i)} = 1 - \sum_{h=1}^i d_h \frac{I(\Phi_l^{(i)} = 1)w_{1h} + I(\Phi_l^{(i)} = 0)w_{0h}}{\tilde{r}_h}, \quad l = 1, \dots, d_i,$$

and

$$A_{ij} = - \sum_{h=1}^i d_h \frac{I(\Phi_{ij} = 1)w_{1h} + I(\Phi_{ij} = 0)w_{0h}}{\tilde{r}_h}, \quad j = 1, \dots, m_i$$

Consider a linear rank statistic

$$v = \sum_{i=1}^k \left(\sum_{l=1}^{d_i} a_l^{(i)} \Phi_l^{(i)} \right) + \sum_{j=1}^{m_i} A_{ij} \Phi_{ij}$$

Equivalence of linear rank statistic & the Gray's statistic

$$\begin{aligned}
 & \quad \quad \quad \vee \\
 &= \sum_{l=1}^{d_1} \left(1 - d_1 \frac{I(\Phi_l^{(1)} = 1)w_{11} + I(\Phi_l^{(1)} = 0)w_{01}}{\tilde{r}_1}\right) \Phi_l^{(1)} - \sum_{j=1}^{m_1} d_1 \frac{I(\Phi_{1j} = 1)w_{11} + I(\Phi_{1j} = 0)w_{01}}{\tilde{r}_1} \Phi_{1j} \\
 &+ \sum_{l=1}^{d_2} \left(1 - d_1 \frac{I(\Phi_l^{(2)} = 1)w_{11} + I(\Phi_l^{(2)} = 0)w_{01}}{\tilde{r}_1} - d_2 \frac{I(\Phi_l^{(2)} = 1)w_{12} + I(\Phi_l^{(2)} = 0)w_{02}}{\tilde{r}_2}\right) \Phi_l^{(2)} \\
 &\quad - \sum_{j=1}^{m_2} \left\{ d_1 \frac{I(\Phi_{2j} = 1)w_{11} + I(\Phi_{2j} = 0)w_{01}}{\tilde{r}_1} + d_2 \frac{I(\Phi_{2j} = 1)w_{12} + I(\Phi_{2j} = 0)w_{02}}{\tilde{r}_2} \right\} \Phi_{2j} \\
 & \quad \quad \quad + \dots + \\
 &+ \sum_{l=1}^{d_k} \left(1 - d_1 \frac{I(\Phi_l^{(k)} = 1)w_{11} + I(\Phi_l^{(k)} = 0)w_{01}}{\tilde{r}_1} - \dots - d_k \frac{I(\Phi_l^{(k)} = 1)w_{1k} + I(\Phi_l^{(k)} = 0)w_{0k}}{\tilde{r}_k}\right) \Phi_l^{(k)} \\
 &\quad - \sum_{j=1}^{m_k} \left\{ d_1 \frac{I(\Phi_{kj} = 1)w_{11} + I(\Phi_{kj} = 0)w_{01}}{\tilde{r}_1} + \dots + d_k \frac{I(\Phi_{kj} = 1)w_{1k} + I(\Phi_{kj} = 0)w_{0k}}{\tilde{r}_k} \right\} \Phi_{kj}
 \end{aligned}$$

Equivalence of linear rank statistic & the Gray's statistic

$$\begin{aligned}
&= \left\{ \sum_{l=1}^{d_1} \Phi_l^{(1)} + \sum_{l=1}^{d_2} \Phi_l^{(2)} + \cdots + \sum_{l=1}^{d_k} \Phi_l^{(k)} \right\} \\
&- d_1 \frac{w_{11}}{\tilde{r}_1} \left\{ \sum_{l=1}^{d_1} \Phi_l^{(1)} + \sum_{l=1}^{d_2} \Phi_l^{(2)} + \cdots + \sum_{l=1}^{d_k} \Phi_l^{(k)} + \sum_{j=1}^{m_1} \Phi_{1j} + \sum_{j=1}^{m_2} \Phi_{2j} + \cdots + \sum_{j=1}^{m_k} \Phi_{kj} \right\} \\
&- d_2 \frac{w_{12}}{\tilde{r}_2} \left\{ \sum_{l=1}^{d_2} \Phi_l^{(2)} + \cdots + \sum_{l=1}^{d_k} \Phi_l^{(k)} + \sum_{j=1}^{m_2} \Phi_{2j} + \cdots + \sum_{j=1}^{m_k} \Phi_{kj} \right\} \\
&\quad \dots \\
&- d_k \frac{w_{1k}}{\tilde{r}_k} \left\{ \sum_{l=1}^{d_k} \Phi_l^{(k)} + \sum_{j=1}^{m_k} \Phi_{kj} \right\} \\
&= \left\{ \sum_{l=1}^{d_1} \Phi_l^{(1)} + \sum_{l=1}^{d_2} \Phi_l^{(2)} + \cdots + \sum_{l=1}^{d_k} \Phi_l^{(k)} \right\} - \left\{ d_1 \frac{\tilde{r}_{11}}{\tilde{r}_1} + d_2 \frac{\tilde{r}_{12}}{\tilde{r}_2} + \cdots + d_k \frac{\tilde{r}_{1k}}{\tilde{r}_k} \right\} \\
&= \sum_{i=1}^k \left(\sum_{l=1}^{d_i} \Phi_l^{(i)} - d_i \frac{\tilde{r}_{1i}}{\tilde{r}_i} \right) \\
&\equiv \tilde{U}_1
\end{aligned}$$

Proposed test approach

- WLOG, $Z_1 < Z_2 < \dots < Z_n$: ordered prognostic factors
- For a given cutpoint μ , $g_i = I(Z_i > \mu)$, i.e., all subjects having Z_i greater than μ belong to group 1, and 0 otherwise
 - As μ , take Z_1, Z_2, \dots , and Z_n , sequentially
- $s_{\mu:1}, s_{\mu:2}, \dots, s_{\mu:n}$: Gray-statistic-type scores of T_1, T_2, \dots, T_n associated with the ranked prognostic factors
- Under H_0 that X & Z are independent and random censoring,
 - $s_{\mu:1}, s_{\mu:2}, \dots, s_{\mu:n}$ are exchangeable
 - $\sum_{i=1}^n s_{\mu:i} = 0$
 - The Noether's condition is satisfied:

$$\frac{\max_{1 \leq i \leq n} |s_{\mu:i}|}{\sum_{i=1}^n s_{\mu:i}^2} \rightarrow 0$$

in probability

Proposed test approach

- Define a process

$$L_n(t) = \frac{\sum_{i=1}^{\lfloor nt \rfloor} S_{\mu:i}}{\sqrt{n-1}\sigma} = \frac{\tilde{U}_1}{\sqrt{n-1}\sigma}, \quad t = 0, 1/n, 2/n, \dots, 1$$

- $\sigma^2 = \frac{1}{n-1} \sum_{i=1}^n S_{\mu:i}^2$
 - $L_n(t)$: re-scaled Gray statistic
- Estimated cutpoint value, $\hat{\mu}$: value of μ which yields the maximum of $|L_n(t)|$
- Significance test
 - Under H_0 , $L_n(t)$ converges in distribution to the Brownian bridge (Billingsley, 1999)
 - For $Q = \sup_{t \in [0,1]} |L_n(t)|$,

$$\text{p-value} = \Pr(Q > q) = 2 \sum_{j=1}^{\infty} (-1)^{j+1} \exp(-2j^2 q^2)$$

Simulations: setup

- Prognostic factor: $Z \sim U(0, 1)$
- True cutpoint: $\mu = 0.2, 0.4, 0.5$
- Effect size: $\beta = 0$ (null); 1, 2, 5(alternative)
- Assumed PH model: $\lambda(t|Z > \mu) = \exp(\beta)\lambda(t|Z \leq \mu)$
- Time to event of interest(X) and time to a competing risk(Y) were generated from the Gumble's bivariate exponential distribution with a degree of dependency, α . Set $\alpha = 0, 0.3$
- Censoring times(C) were independently generated with time to event of interest or time to a competing risk from an exponential distribution with hazard rate of k
 - Censoring fraction: $p = 0, 0.3$
 - $k : P(C < X \wedge Y) = p$, i.e.,

$$k : \mu k \left\{ \frac{1+\alpha}{2+k} - \frac{2\alpha}{3+k} + \frac{\alpha}{4+k} \right\} + (1-\mu) k \left\{ \frac{1+\alpha}{\theta+1+k} - \frac{\alpha}{\theta+2+k} - \frac{\alpha}{2\theta+1+k} + \frac{\alpha}{2\theta+2+k} \right\}$$

$$= p$$

Data generation procedure

- Step 0: Fix μ , $\theta (= e^\beta)$, α , and p
- Step 1: generate a random variate z from $U(0, 1)$
- Step 2: generate a random variate u_1 from $U(0, 1)$
 - If $z \geq \mu$, $x = -\ln(1 - u_1)/\theta$
 - If $z < \mu$, $x = -\ln(1 - u_1)$
- Step 3: generate a random variate u_2 from $U(0, 1)$

$$y = -\ln\left[\frac{w(x)-1+\sqrt{(w(x)-1)^2+4w(x)(1-u_2)}}{2w(x)}\right]$$

- $w(x) = d_1(x)c_1(x) + d_2(x)c_2(x)$ with $c_1(x) = \alpha(2e^{-x} - 1)$,
 $c_2(x) = \alpha(2e^{-\theta x} - 1)$, $d_1(x) = \frac{\mu f_1(x)}{\mu f_1(x) + (1-\mu)\theta f_1(\theta x)}$, and
 $d_2(x) = \frac{(1-\mu)\theta f_1(\theta x)}{\mu f_1(x) + (1-\mu)\theta f_1(\theta x)}$
- f_1 : pdf of $\exp(1)$

Data generation procedure

- Step 4: generate a random variate c from $\exp(k)$
- Step 5: define

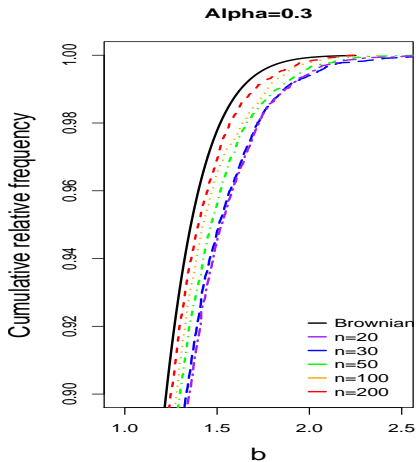
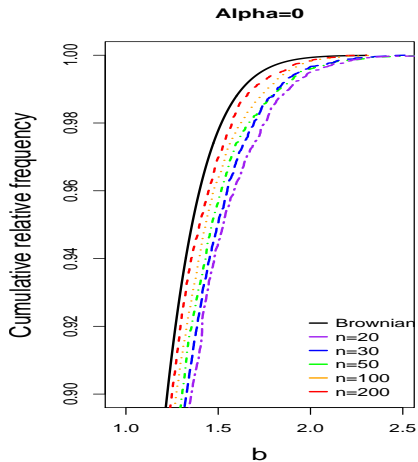
$$t = (x \wedge y) \wedge c$$

and

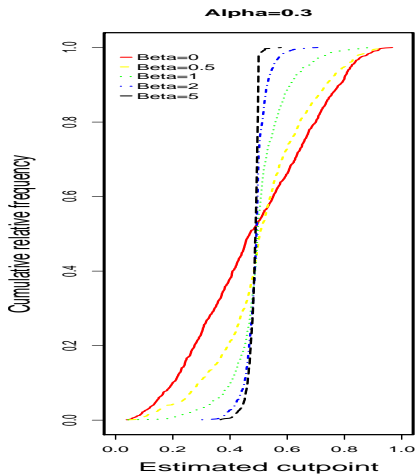
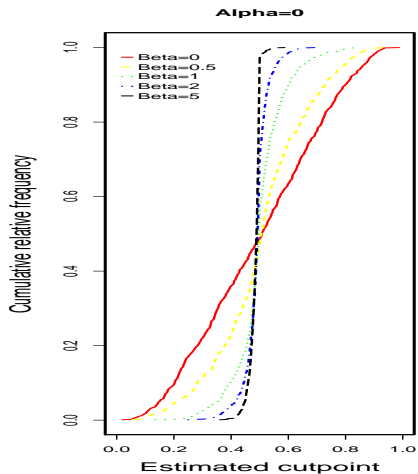
$$\delta = 1 \times I(x < y \wedge c) + 2 \times I(y < x \wedge c),$$

i.e., event occurred ($\delta = 1$), competing risk(s) occurred ($\delta = 2$),
censored ($\delta = 0$)

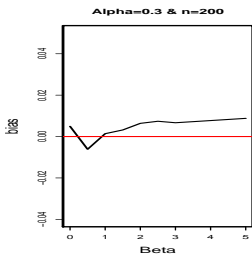
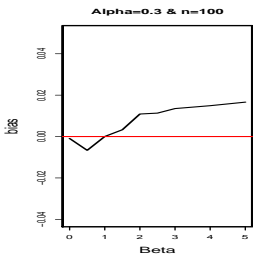
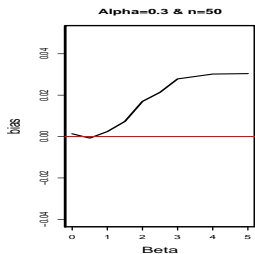
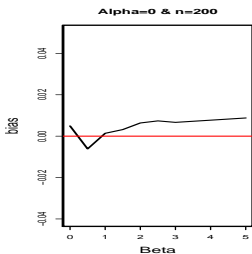
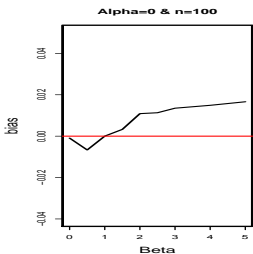
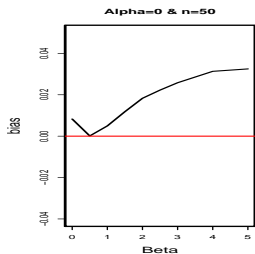
Upper part of the distribution of the extreme value of the standardized process when $\mu = 0.5$



Distribution of the cutpoint estimator with $n = 100$ and $\mu = 0.5$



Average bias(=true-estimated) with $\mu = 0.5$



Simulation results: when $\alpha = 0$

Table : Empirical bias(bias) and standard deviation(SD) of the estimated cutpoint and power(power) of the proposed test under no censoring

μ	β	CR	$n = 50$			$n = 100$			$n = 200$		
			bias	SD	power	bias	SD	power	bias	SD	power
0.2	0	50	-0.292	0.209	0.040	-0.297	0.218	0.038	-0.296	0.218	0.044
	1	31	-0.189	0.206	0.229	-0.139	0.176	0.484	-0.091	0.136	0.817
	2	20	-0.057	0.118	0.754	-0.036	0.074	0.984	-0.022	0.046	1.000
	5	10	-0.001	0.043	0.989	0.001	0.020	1.000	0.000	0.012	1.000
0.4	0	50	-0.096	0.213	0.047	-0.103	0.211	0.045	-0.091	0.221	0.037
	1	36	-0.045	0.137	0.538	-0.034	0.108	0.816	-0.023	0.073	0.979
	2	27	-0.001	0.064	0.992	-0.001	0.044	1.000	-0.001	0.028	1.000
	5	20	0.023	0.026	1.000	0.012	0.016	1.000	0.006	0.007	1.000
0.5	0	50	0.019	0.212	0.040	0.011	0.225	0.046	0.007	0.213	0.042
	1	38	0.011	0.123	0.491	0.000	0.090	0.814	0.000	0.053	0.989
	2	31	0.022	0.054	0.986	0.010	0.030	1.000	0.004	0.016	1.000
	5	25	0.029	0.030	1.000	0.014	0.016	1.000	0.007	0.007	1.000

Simulation results: when $\alpha = 0.3$

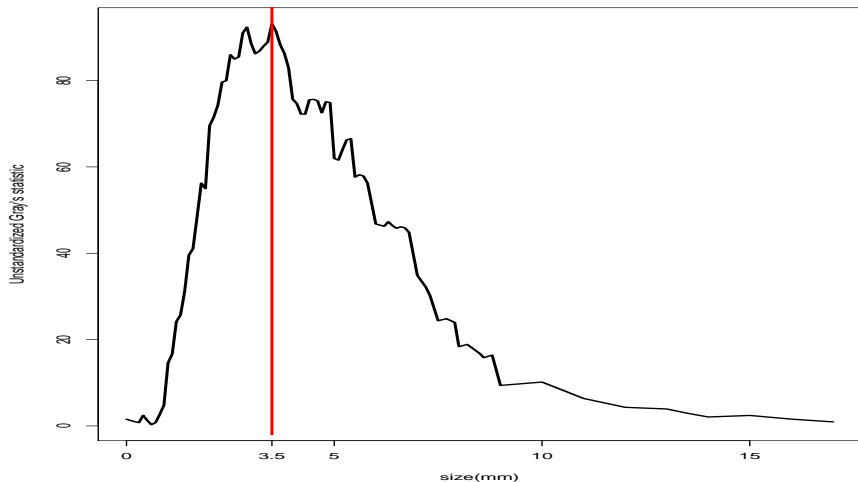
Table : Empirical bias and standard deviation of the estimated cut-off value and power of the proposed test under no censoring

μ	β	CR	$n = 50$			$n = 100$			$n = 200$		
			bias	SD	power	bias	SD	power	bias	SD	power
0.2	0	50	-0.287	0.219	0.031	-0.292	0.217	0.040	-0.297	0.216	0.032
	1	31	-0.191	0.197	0.319	-0.151	0.175	0.530	-0.094	0.134	0.796
	2	19	-0.076	0.138	0.798	-0.059	0.098	0.981	-0.040	0.072	1.000
	5	11	-0.001	0.041	0.998	-0.002	0.029	1.000	-0.001	0.014	1.000
0.4	0	50	-0.087	0.219	0.036	-0.090	0.215	0.033	-0.102	0.211	0.041
	1	35	-0.048	0.143	0.523	-0.031	0.104	0.802	-0.020	0.071	0.972
	2	27	-0.005	0.070	0.984	-0.001	0.046	1.000	-0.002	0.032	1.000
	5	20	0.023	0.026	1.000	0.013	0.014	1.000	0.006	0.008	1.000
0.5	0	50	0.002	0.217	0.036	0.016	0.217	0.055	0.004	0.216	0.040
	1	38	-0.005	0.130	0.552	-0.003	0.091	0.810	-0.001	0.063	0.985
	2	30	0.018	0.061	0.985	0.009	0.038	1.000	0.006	0.025	1.000
	5	25	0.032	0.032	1.000	0.018	0.020	1.000	0.008	0.009	1.000

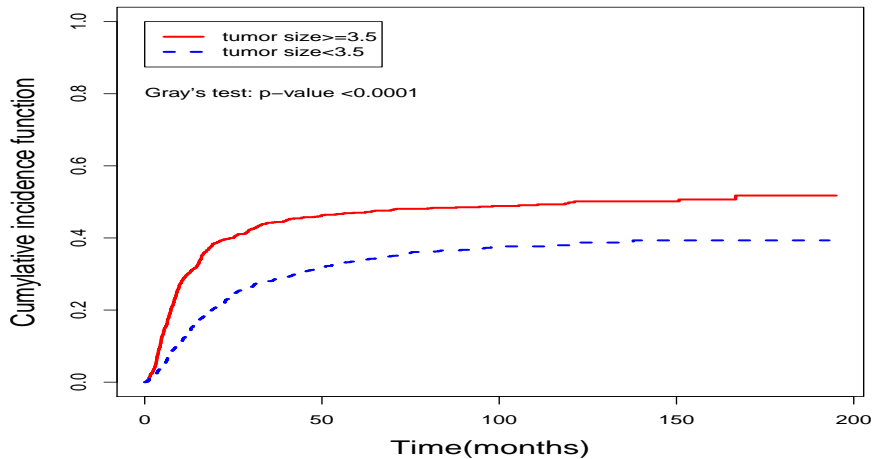
Lung cancer data(Revisited)

- The 1965 lung cancer patients received a tumor removal surgery at Seoul Samsung Hospital in Korea from September 1994 to December 2005
- Event of interest: relapse, competing risk(s): death, prognostic factor: tumor size at surgery
- Relapsed(43%), death(18%), censored(39%)
- Tumor size: ranged over (0,19)

Determining the optimal cutpoint



Cumulative incidence functions by tumor size



Discussion

- Represent the Gray's statistic as a linear rank statistic in competing risks model
- Propose a procedure for identifying the cutpoint and significance test
- Simulations showed that the null distribution of our test statistic is getting close to the theoretical distribution as the sample size increases, and the proposed test satisfied the nominal level and the power of our proposed test increases as the effect size increases.
- We found that the optimal cutpoint of tumor size is 3.5mm, which is in accordance with physicians' experience
- Limitations: extensive simulations are required under the moderate to heavy censoring to investigate the small-sample properties of our test

Thank you!