



Published in final edited form as:

*Biometrics*. 2018 June ; 74(2): 694–702. doi:10.1111/biom.12770.

## Efficiency of Two Sample Tests via the Restricted Mean Survival Time for Analyzing Event Time Observations

Lu Tian<sup>1</sup>, Haoda Fu<sup>2</sup>, Stephen J. Ruberg<sup>2</sup>, Hajime Uno<sup>3</sup>, and LJ Wei<sup>4</sup>

<sup>1</sup>Department of Biomedical Data Science, Stanford University, CA 94305, U.S.A

<sup>2</sup>Eli Lilly and Company, Indianapolis, IN 46285, U.S.A

<sup>3</sup>Dana-Faber/Harvard Cancer Institute, Boston, MA 02215, U.S.A

<sup>4</sup>Department of Biostatistics, Harvard University, Boston, MA 02115, U.S.A

### Summary

In comparing two treatments with the event time observations, the hazard ratio (HR) estimate is routinely used to quantify the treatment difference. However, this model dependent estimate may be difficult to interpret clinically especially when the proportional hazards (PH) assumption is violated. An alternative estimation procedure for treatment efficacy based on the restricted means survival time or  $t$ -year mean survival time ( $t$ -MST) has been discussed extensively in the statistical and clinical literature. On the other hand, a statistical test via the HR or its asymptotically equivalent counterpart, the logrank test, is asymptotically distribution-free. In this paper, we assess the relative efficiency of the hazard ratio and  $t$ -MST tests with respect to the statistical power under various PH and non-PH models theoretically and empirically. When the PH assumption is valid, the  $t$ -MST test performs almost as well as the HR test. For non-PH models, the  $t$ -MST test can substantially outperform its HR counterpart. On the other hand, the HR test can be powerful when the true difference of two survival functions is quite large at end but not the beginning of the study. Unfortunately, for this case, the HR estimate may not have a simple clinical interpretation for the treatment effect due to the violation of the PH assumption.

### Keywords

Asymptotic relative efficiency; Hazard ratio; Proportional hazards model; Survival analysis;  $t$ -year mean survival time

## 1. Introduction

In a randomized, comparative clinical trial with an event time as the study end point, treatment difference is often summarized by the hazard ratio (HR) by assuming a proportional hazards (PH) model, namely, the ratios of two hazard functions between two treatment groups, labeled as 1 and 0, are approximately constant over time (Cox, 1972). Under the PH assumption, the HR can be consistently estimated by maximizing a partial

likelihood function from the Cox model. However, a HR of, for instance, 0.7 for treatment 1 versus treatment 0 cannot be interpreted as a 30% event rate reduction in favor of treatment 1 since the hazard is not a simple probability of event occurrence. The clinical interpretation of a HR depends on the temporal profile of the hazard function for the reference group. For example, it is not clear whether a HR of 0.7 alone is clinically meaningful: a 30% reduction of hazard in treatment group 1 may not be clinically important for a “low” hazard function in treatment group 0, while quite meaningful for a “high” hazard function in the treatment group 0. However, there is a lack of a simple nonparametric estimator for the hazard function as the Kaplan-Meier (KM) estimator for the survival function. Oftentimes the HR is interpreted as the inverse ratio of the median event times in practice assuming an exponential distribution for the survival outcome and thus provides an informative summary for the treatment effect when coupled with the median survival time in group 0. However, such an interpretation may be problematic when the exponential assumption does not hold. Furthermore, when the PH assumption is violated, the HR estimator from the Cox model converges to a parameter which is difficult to interpret (Kalbeisch and Prentice, 1981; Lin and Wei, 1989). For this case, one often considers the estimated HR as an approximation to a weighted average of the HR’s over time. Unfortunately, the weights depend on the censoring distributions. This adds another complexity of translating HR for effective clinical decision making. Other model-based summary measures for the group difference have similar issues as the HR (Wei, 1992).

There are several alternatives to summarize a survival distribution. For example, the median survival time, the  $t$ -year event rate and the  $t$ -year mean survival time,  $t$ -MST. The  $t$ -year mean survival time is also coined as the restricted mean survival time. In this paper we are interested in studying properties of using a summary measure as a group contrast based on the  $t$ -MST. The inference procedures for  $t$ -MST and the function thereof have been studied extensively, for example, by Karrison (1987); Zucker (1998); Royston and Parmar (2011); Zhao et al. (2012); Tian et al. (2014) and Uno et al. (2014). The  $t$ -MST has a clear physical and clinical interpretation. For example, an observed 2.5 year of the 3-year mean survival time indicates if a patient is followed up to 3 years, on average, he or she would survive 2.5 years. The  $t$ -MST can be readily estimated via the area under the corresponding Kaplan-Meier (KM) curve up to  $t$  year. Contrary to the model-dependent nature of the HR, the treatment effect can be quantified by, for example, the difference in  $t$ -MST between two treatment groups, which is purely nonparametric without requiring any model assumption. This group difference measure with a reference value of  $t$ -MST from the control arm in a comparative study is more informative than the HR.

Since the HR estimate is routinely used for making inference about the treatment effect, it would be interesting to compare it with the  $t$ -MST based inference procedure. However, since these two estimation procedures empirically quantify different parameters, it seems difficult, if not impossible, to study their relative merits from an estimation point of view. On the other hand, under the testing hypothesis paradigm, the HR based test is asymptotically distribution-free. Therefore, it is possible to compare the  $t$ -MST and HR based tests with respect to, for example, their conventional statistical power profiles. It is interesting to note that recently Trinquart et al. (2016) utilized the data from a large number of clinical trials to

show empirically these two types of tests are quite concordant in terms of the conventional statistical significance interpretation.

In this paper, we formally compare these two types of tests under a more theoretical setting. To be specific, we first present the testing procedures based on HR and  $t$ -MST in Section 2. The asymptotic relative efficiency (ARE) of these two tests is derived in Section 3. In Section 4, the performance of the two tests are studied theoretically and compared empirically via extensive numerical study. We also use the data from a recent clinical trial to illustrate our findings in the same section. In general, we find that the test based on  $t$ -MST performs well compared with its HR counterpart when the PH assumption is valid. For certain non-PH models, for instance, when the two survival functions are quite large at the end but not the beginning of the study followup time, the HR test is quite powerful. Unfortunately, the corresponding HR estimate is difficult to interpret clinically. With this additional finding, the robust, easily interpretable  $t$ -MST based inference procedure provides a useful alternative tool to the HR based counterpart in survival analysis.

## 2. The HR and $t$ -MST based tests for the equivalence of two survival functions

Let  $T_1$  and  $T_0$  be the event times in treatment groups 1 and 0, respectively. The group difference in  $t$ -MST up to the time point  $t$  is defined as

$$D = E(T_1 \wedge t) - E(T_0 \wedge t),$$

which is the area between two survival curves over the time interval  $[0, t]$ :

$$\int_0^t \{S_1(u) - S_0(u)\} du,$$

where  $S_j(\cdot)$  is the survival function of the failure time  $T_j$  in arm  $j$ ,  $j = 0, 1$ . In order to make  $D$  identifiable based on observed data,  $\pi_j(t) = \text{pr}(T_j \wedge C_j > t)$  needs to be bounded away from zero, where  $C_j$  denotes the censoring time in treatment group  $j$ ,  $j = 0, 1$ , i.e.,  $\min\{\pi_0(t), \pi_1(t)\} > 0$ . In practice,  $t$  can be chosen, for example, as the minimum of the 95th percentiles of observed  $T_j \wedge C_j$  from two treatment groups. The difference in  $t$ -MST can then be estimated by

$$\hat{D} = \int_0^t \{\hat{S}_1(u) - \hat{S}_0(u)\} du,$$

where  $\hat{S}_j(t)$  is the KM estimator for  $S_j(t)$ . As the sample size  $n \rightarrow \infty$ ,  $\sqrt{n}(\hat{D} - D)$  converges weakly to a mean zero Gaussian distribution whose variance can be consistently estimated by

$$\hat{\sigma}_D^2 = \int_0^t \left\{ \int_v^t \hat{S}_0(u) du \right\}^2 \frac{d\hat{\Lambda}_0(v)}{p_0 \hat{\pi}_0(v)} + \int_0^t \left\{ \int_v^t \hat{S}_1(u) du \right\}^2 \frac{d\hat{\Lambda}_1(v)}{p_1 \hat{\pi}_1(v)}, \quad (1)$$

where  $p_j$  is the proportion of patients randomized into group  $j$ ,  $\hat{\pi}_j(v)$  is the empirical counterpart of  $\pi_j(v)$  and  $\hat{\Lambda}_j(t)$  is the Nelson-Aalen estimator for the cumulative hazard function in group  $j$ ,  $j = 0, 1$  (Pepe and Fleming, 1989, 1991; Zhao et al., 2012). Under the null hypothesis that there is no difference between two survival functions, the distribution of  $\sqrt{n}\hat{D}/\hat{\sigma}_D$  is approximately the standard normal for large  $n$ .

Similarly, let  $\hat{\theta}$  be the estimator of  $\theta = \log(\text{HR})$  by maximizing the partial likelihood function over the time interval  $[0, t_H]$ . Under the PH model, as  $n \rightarrow \infty$ ,  $\sqrt{n}(\hat{\theta} - \theta)$  converges weakly to a mean zero Gaussian distribution, whose variance can be consistently estimated by

$$\hat{\sigma}_\theta^2 = \left\{ \int_0^{t_H} \frac{e^{\hat{\theta}} p_0 p_1 \hat{\pi}_0(u) \hat{\pi}_1(u)}{p_0 \hat{\pi}_0(u) + e^{\hat{\theta}} p_1 \hat{\pi}_1(u)} d\tilde{\Lambda}_0(u) \right\}^{-1}, \quad (2)$$

where  $\tilde{\Lambda}_0(u)$  is the Breslow estimator for the cumulative hazard function at the treatment group 0 (Cox, 1972; Fleming and Harrington, 2011). Under the null hypothesis,  $\theta = 0$  and the distribution of  $\sqrt{n}\hat{\theta}/\hat{\sigma}_\theta$  is approximately standard normal for large  $n$ . Therefore, tests based on both  $\hat{D}$  and  $\hat{\theta}$  are asymptotical valid in retaining the appropriate type I error rate.

Note that similar to  $t$ -MST, in theory the large sample normal approximation to the distribution of HR estimator is only valid with observations within a finite time interval, say  $[0, t_H]$ , where  $\max\{\pi_0(t_H), \pi_1(t_H)\} > 0$  (p289–290, Chapter 8, Fleming and Harrington (2011)). There is a misconception that the  $\exp(\hat{\theta})$  approximates the HR at all time points. In practice, the interpretation of this estimated HR should be restricted to the finite time interval  $[0, t_H]$ , where  $t_H$  is smaller than the maximum the followup time. Compared with  $t$  in  $t$ -MST,  $t_H$  and  $t$  can be identical under the condition that  $\sup\{u \mid \pi_1(u) > 0\} = \sup\{u \mid \pi_0(u) > 0\}$ . Otherwise, without loss of generality assume that  $\sup\{u \mid \pi_0(u) > 0\} < \sup\{u \mid \pi_1(u) > 0\}$ . In such a case, one may choose  $t_H = \sup\{u \mid \pi_0(u) > 0\}$ , which is greater than  $t = \sup\{u \mid \pi_0(u) > \varepsilon\}$  for some positive  $\varepsilon$ . However, in theory, the difference between  $t$  and  $t_H$  can be made small by choosing  $\varepsilon$  close to zero.

### 3. Asymptotical efficiency of the test based on $\hat{D}$ and $\hat{\theta}$

In this section, we derive the asymptotic efficiencies of the tests based on  $\hat{D}$  and  $\hat{\theta}$  under a sequence of contiguous alternatives (Hoeffding and Rosenblatt, 1955; Hodges and Lehmann, 1956). Specifically, for the current case, these alternatives are defined as

$$\log \left\{ \frac{\lambda_{1n}(u)}{\lambda_0(u)} \right\} = \frac{\alpha(u)}{\sqrt{n}},$$

where  $\lambda_0(t)$  and  $\lambda_{1n}(t)$  are the hazard functions of the failure time  $T_0$  from the treatment 0 and  $T_{1n}$  from the treatment 1, whose distribution depends on the sample size  $n$ , respectively (Lagakos, 1988; Slud, 1991). Here,  $\alpha(\cdot)$  is a deterministic function over time providing a specific alternative hypothesis of interest. For any given  $\alpha(\cdot)$ , we can derive the asymptotic efficiency of the test under the mild regularity conditions given in p230–232, Chapter 6, Fleming and Harrington (2011).

For study size  $n$ , it follows from similar arguments in Schoenfeld (1981) and Harrington and Fleming (1982) that under the above contiguous alternatives  $\sqrt{n}\hat{D}$  is asymptotic normal with mean

$$\int_0^t \left\{ \int_0^v \alpha(u)\lambda_0(u)du \right\} S_0(v)dv$$

and variance

$$\sigma_D^2 = \int_0^t \left\{ \int_v^t S_0(u)du \right\} \frac{2\lambda_0(v)}{w(v)} dv, \quad (3)$$

which is the limit of  $\hat{\sigma}_D^2$  defined in (1), where

$$w(v) = \frac{p_0 p_1 \pi_0(v) \pi_1(v)}{p_0 \pi_0(v) + p_1 \pi_1(v)}.$$

Thus, the asymptotic efficiency of the test based  $\hat{D}$  is

$$\xi_1 = \frac{\left[ \int_0^t \left\{ \int_v^t S_0(u)du \right\} \alpha(v)\lambda_0(v)dv \right]^2}{\int_0^t \left\{ \int_v^t S_0(u)du \right\}^2 w(v)^{-1} \lambda_0(v)dv}.$$

Similarly, it follows from Lin and Wei (1989) that under above contiguous alternatives the test statistic  $\sqrt{n}\hat{\theta}$  is also asymptotic normal with mean

$$\frac{\int_0^t H w(u)\alpha(u)\lambda_0(u)du}{\int_0^t w(u)\lambda_0(u)du}$$

and variance

$$\sigma_{\theta}^2 = \left\{ \int_0^{tH} w(u)\lambda_0(u)du \right\}^{-1}, \quad (4)$$

which is the limit of  $\hat{\sigma}_{\theta}^2$  in (2). The asymptotic efficiency for this test is

$$\xi_2 = \frac{\left\{ \int_0^{tH} w(u)\alpha(u)\lambda_0(u)du \right\}^2}{\int_0^{tH} w(u)\lambda_0(u)du}.$$

It follows that the ARE of the  $t$ -MST based test relative to the HR-based test is  $\xi_1/\xi_2$ , which can be interpreted approximately as the inverse of the ratio of the sample sizes needed for two tests to have the same power to detect the alternative  $\lambda_{1n}(u) = \exp\{\alpha(u)/\sqrt{n}\}\lambda_0(u)$ .

Note that the standard logrank test, a score test based on the Cox partial likelihood function for testing the equivalence of two survival functions, is asymptotically equivalent to the Wald-type of test based on the HR and therefore  $\xi_1/\xi_2$  is also the ARE of the  $t$ -MST based test relative to the logrank test.

Since ARE is approximately a ratio of the sample size of two tests to have similar power for testing the same alternative hypothesis, this connection provides a clear practical interpretation of ARE (Schoenfeld and Richter, 1982; Zhang and Quan, 2009). Specifically, assuming that the alternative hypothesis consists of two fixed unequal hazard functions, say,  $\lambda_0(u)$  and  $\lambda_1(u)$ , the sample size needed for the  $t$ -MST based test can be approximated by  $(z_{1-\alpha/2} + z_{1-\beta})^2/\xi_1$ , if this alternative is “close” to the null hypothesis, where  $z_q$  is  $q$ th quantile of the standard normal and  $\alpha$  and  $\beta$  are the Type I and II error rates, respectively. More generally, a more accurate sample size estimator for  $100(1 - \beta)\%$  power at the two-sided significance level of  $\alpha$  is

$$(z_{1-\alpha/2} + z_{1-\beta})^2 \frac{\int_0^t \left[ \int_v^t S_0(u)du \right]^2 \frac{\lambda_0(v)}{p_0\pi_0(v)} + \left[ \int_v^t S_1(u)du \right]^2 \frac{\lambda_1(v)}{p_1\pi_1(v)} dv}{\left[ \int_0^t \{S_1(u) - S_0(u)\}du \right]^2}, \quad (5)$$

for any pair of  $\{\lambda_0(u), \lambda_1(u)\}$ . Similarly, the corresponding sample size estimate analogy for HR based test is

$$(z_{1-\alpha/2} + z_{1-\beta})^2 \frac{\int_0^t w(u)^2 \left\{ \frac{\lambda_0(u)}{p_0\pi_0(u)} + \frac{\lambda_1(u)}{p_1\pi_1(u)} \right\} du}{\left[ \int_0^t w(u)\{\lambda_1(u) - \lambda_0(u)\}du \right]^2}. \quad (6)$$

When  $\lambda_1(u)$  for the treated arm is close to  $\lambda_0(u)$  for the control arm, the ratio of the two sample size estimates is approximately equal to the inverse of ARE. The ARE, coupled with these two sample size estimates, may provide a meaningful comparison of these two tests. However, when  $\lambda_1(u)$  is quite different from  $\lambda_0(u)$ , ARE may not be a good approximation to such a ratio of the sample sizes. In this case, one may consider the empirical relative efficiency (ERE),  $E_D^2/E_\theta^2$ , where

$$E_D = \frac{E(\hat{D})}{se(\hat{D})} \text{ and } E_\theta = \frac{E(\hat{\theta})}{se(\hat{\theta})}$$

are the effect sizes standardized by their standard errors for tests based on  $\hat{D}$  and  $\hat{\theta}$ , respectively. While ERE in general can only be approximated numerically, a close-form expression for ARE can often be obtained via a series of first order approximations under the contiguous alternative. For example, the approximation

$$\begin{aligned} E(\hat{D}) &\approx \int_0^t \{S_1(u) - S_0(u)\} du = \int_0^t S_0(v) \left[ e^{-\int_0^v \{\lambda_1(u) - \lambda_0(u)\} du} - 1 \right] dv \\ &\approx \int_0^t \left[ \int_0^v \lambda_0(u) \log\{\lambda_0(u)/\lambda_1(u)\} du \right] S_0(v) dv, \end{aligned}$$

may be used to compute the asymptotic efficiency of the test based on  $t$ -MST. While the first approximation above is fairly accurate in general, the second approximation is sensitive to the distance between  $\lambda_1(u)$  and  $\lambda_0(u)$ . Similar approximation has been employed to derive  $\sigma_D^2$  in (3). Consequently, the asymptotic efficiency  $\xi_1$  may not be a good approximation to the ratio  $\{E(\hat{D})\}^2/\text{var}(\hat{D})$ , which ultimately dictates the actual power of the test based on  $t$ -MST, especially when  $\lambda_1(u)$  and  $\lambda_0(u)$  are markedly different at some time points. On the other hand, for any given sample size, we may simulate a large number of  $\hat{D}$ s according to the postulated survival and censoring distributions under the alternative and directly approximate  $E(\hat{D})$  and  $\text{var}(\hat{D})$  by the empirical average and variance of simulated  $\hat{D}$ s, respectively, without relying on any assumption on the difference between  $\lambda_1(u)$  and  $\lambda_0(u)$ . This approach is even more appealing in estimating the expectation and variance of the HR estimator, which is the solution of a complex nonlinear estimating equation and does not have an exact close-form expression. In summary, although estimating ERE is more computational intensive, the power calculation based on ERE is more robust and the approximations to (5) and (6) can be used to evaluate the relative merits of two tests under any alternative. An example is given in the next section to illustrate this point.

#### 4. Example and Numerical Study

In this section, we will evaluate the empirical performance of these two tests. To gain insights for the numerical results to be presented, we first note that under the contiguous alternative, the  $t$ -MST based test is asymptotically equivalent to the test based on

$$\int_0^t \frac{\int_v^t S_0(u) du}{\int_0^t S_0(u) du} d\left\{ \widehat{\Lambda}_1(v) - \widehat{\Lambda}_0(v) \right\},$$

while the HR-based test is equivalent to the test based on

$$\int_0^{t_H} S_{C(v)} S_0(v) d\left\{ \widehat{\Lambda}_1(v) - \widehat{\Lambda}_0(v) \right\}.$$

Here we assume that the survival functions of the censoring distributions at two groups are identical and denoted by  $S_C(\cdot)$ . If  $t = t_H$ , the difference between these two tests thus only depends on the choice of weight function. To be specific, let

$$w_D(v) = \frac{\int_v^t S_0(u) du}{\int_0^t S_0(v) dv} \quad \text{and} \quad w_\theta(v) = S_C(v) S_0(v)$$

be the weight functions of the  $t$ -MST and HR based tests, respectively.  $w_D(v)$  is independent of the censoring distribution and monotone decreasing from 1 to 0, while  $w_\theta(v)$  depends on the censoring distribution and monotone decreasing from 1 to  $S_C(t) S_0(t) > 0$ . In Figure 2, we plot the weight functions for both tests with following survival functions for the failure and censoring time distributions:

$$S_0(v) = \exp(-0.016v^{0.826}) \quad \text{and} \quad S_C(v) = e^{-\lambda_c v \frac{\{\tau_u - \max(v, \tau_l)\}}{(\tau_u - \tau_l)}} I(v \leq \tau_u),$$

where  $\lambda_c = 0.06\%$  or  $2.5\%$  and  $[\tau_l, \tau_u] = [24, 43]$ . Those model parameters are selected to mimic the ECOG-ACRIN study described in the next subsection. It is clear that  $w_D(v)$  assigns less weights to later difference in hazard function compared with  $w_\theta(v)$  when the censoring is light. In such a case, the relative performance of  $t$ -MST and HR based tests should depend on the temporal profile of the difference between two hazard functions. For example, when the treatment difference is anticipated at later study time points, the  $t$ -MST based test tends to be less powerful than the HR based test. On the other hand, when the drop-out rate is high,  $w_D(v)$  can be very similar to  $w_\theta(v)$ , suggesting similar performance of these two tests. In the next two subsections, we will present detailed results from numerical study, which confirms the aforementioned expectation.

### 4.1 Real Data Example

We first use a study recently conducted by the ECOG-ACRIN Cancer Research Group to compare low- and high-dose dexamethasone for treating newly diagnosed multiple myeloma (Rajkumar et al., 2010) to illustrate how to interpret the ARE ( $t$ -MST vs. HR based tests). This study randomized 222 patients to the low-dose group and 223 patients to the high-dose group. Figure 1a shows the resulting KM curves of overall survival by treatment groups. Visually it seems that the patients in the low-dose group survived longer than those in the



high dose group. The p-values from the tests based on HR and  $t$ -MST are 0.47 and 0.04, respectively. The discrepancy is likely due to the presence of crossing hazards. In this study, the nonparametric KM curves are fairly similar to their parametric counterparts based on a Weibull model (Figure 1b). Now, suppose that we are interested in designing a new study using the results from this ECOG-ACRIN study. That is, assume that the observed pattern of the two observed hazard functions from this study is the alternative hypothesis for a new study. The question is how the  $t$ -MST test would perform compared with the HR based test with respect to the ARE and how to interpret this ratio measuring the relative merit of these two tests. To this end, we assume that the event time follows Weibull distributions as shown in Figure 1b. Furthermore, we assume this new study would have similar patient's accrual and follow up time patterns, that is the entire study duration is about 43 months with a 19-month enrollment period. Furthermore, we let the time of loss of follow-up follows an exponential distribution with an annual drop-off rate of 1% for both arms without accounting for the administrative censoring due to the staggered entry of study patients. We also let  $t = 40$  and  $t_H = 43$  months for defining the follow-up period of interest for  $t$ -MST and HR, respectively. Under these assumptions, the required sample size for 80% power at the significance level of 0.05 is  $n_\theta = 2532$  per arm for the HR-based test and  $n_D = 564$  per arm for the  $t$ -MST based test. The ARE is 2.23, which also strongly favors  $t$ -MST based test but quantitatively different  $n_\theta/n_D$ . In this case, we also can use simulation to directly estimate the ERE, the finite sample analogy of ARE. To this end, we simulate 5000 sets of data consisting of 1000 patients each under the assumed alternatives and obtain the corresponding  $\hat{D}$  and  $\hat{\theta}$  from each simulated data set. We then approximate the expectation and variance of  $\hat{D}$  (and  $\hat{\theta}$ ) by their empirical counterparts. The resulting ERE is 4.41, which is pretty close to  $n_\theta/n_D$ . It suggests that the "alternative" of interest is not adequately close to null and ARE may not accurately reflect the ratio of the needed sample sizes of the two tests in this case. With the above sample size estimates, one still can assess the practical gain from the  $t$ -MST test over the HR test. We have performed the additional simulation studies to empirically estimate the true power of HR and  $t$ -MST based tests with the estimated sample sizes. The empirical powers based on 5000 simulations are 80.1% and 80.9% for  $t$ -MST and HR based tests, respectively, which confirms the adequacy of the estimated sample sizes based on (5) and (6). The required sample size is much greater than that in actual study, since more patients are needed to ensure 80% chance of obtaining a p-value smaller than 0.05.

Now, we use the above setting, but modify the underlying Weibull distribution of the low dose group to follow the PH assumption with a true HR of 0.70 against the high-dose arm. The survival curves of this alternative are plotted in Figure 1c. With this specific PH alternative, the required sample sizes are  $n_D = 641$  for  $t$ -MST based test and  $n_\theta = 591$  for HR based test. Both ARE and ERE are 0.96 and 0.91, respectively, which are pretty close to  $n_\theta/n_D$ . This example suggests that even when the PH assumption is valid, the  $t$ -MST test is essentially as powerful as the HR test.

## 4.2 Simulation Study

In this subsection, we further explore extensively the relative merits between the two tests under various scenarios. We focus on the settings, where the survival functions are ordered

without crossing during the followup and we can claim patients from one arm live longer than those from another during the time interval of interest without any ambiguity. When two KM curves cross during the study followup period, one may explore whether there is a subgroup of patients who may not benefit from the new therapy based on the patients baseline covariates. Note that the hazard functions may still cross during the followup even if the survival curves don't. Now, for our study, firstly, we consider the case where  $t_H = t$  and the PH assumption holds true, i.e.,  $\alpha(u) = 1$ . For this case, the ARE

$$\leq \frac{[\int_0^t \int_v^t S_0(u) du \lambda_0(v) dv]^2}{\int_0^t \int_v^t S_0(u) du \int_0^t w(v)^{-1} \lambda_0(v) dv \times \int_0^t w(u) \lambda_0(u) du},$$

which is always  $\geq 1$  by Cauchy inequality, suggesting that the HR-based test is more powerful under the PH assumption. To quantify the efficiency loss of the  $t$ -MST based test under practical settings, we assume that the study has a recruitment period of 19 months and additional follow-up of 24 months after the last patient entered the trial as in the aforementioned ECOG-ACRIN study. Specifically, the censoring time is assumed to be the minimum of  $C_L \sim \text{EXP}(\lambda_c)$  and  $C_A \sim U(\tau_L, \tau_U)$ , reflecting the loss of follow-up and administrative censoring due to staggered entry, respectively. It follows that the survival function for the censoring time is

$$S_C(u) = e^{-\lambda_c u} \frac{(\tau_U - \max(u, \tau_L))}{(\tau_U - \tau_L)} I(u \leq \tau_U),$$

where  $[\tau_L, \tau_U] = [24, 43]$ . Here, we let  $t = 40$  and  $t_H = 42.5$  months. Furthermore, we set  $\lambda_c = 0.06\%$  matching the annual loss of follow-up rate observed in the ECOG-ACRIN study and  $\lambda_c = 2.5\%$ , an annual loss of follow-up rate of 26%, to represent heavier censoring caused by, for example, drop-out. Lastly, we also assume that the survival time in the treatment group 0 follows a Weibull distribution with the shape parameter being the maximum likelihood estimator in the high-dose group of the ECOG-ACRIN study, that is, having a decreasing hazard function  $\lambda_0(u) = a_0 u^{-0.174}$ . With this setup, ARE is determined by  $a_0$ , the scale parameter of the Weibull distribution in the treatment group 0. We adjusted the scale parameter so that  $S_0(t) = 0.1, \dots, 0.9$ . The detailed results on ARE are summarized in Table 1. Under the PH assumption, the HR-based test is slightly more powerful as anticipated. Furthermore, when the cumulative event rate is relatively high, for example,  $\text{pr}(T_0 < t) = 0.5$ , as for studies of serious diseases with high event rates, the efficiency loss of  $t$ -MST based test relative to the HR based optimal test is almost negligible.

To confirm these theoretical AREs, we have also performed numerical simulations to obtain the ERE by estimating  $E_D$  and  $E_R$  based on results from 5000 simulated data sets. The true HR for two treatment groups (treatment group 1 vs. 0) is 0.70. The resulting EREs are all very close to their asymptotical counterparts (Table 1). For example, when  $\lambda_c = 0.06\%$  and  $S_0(t) = 0.60$ , the ARE and ERE are 0.98 and 0.96, respectively.

Next, we consider three non-PH settings:

1. HR is monotone increasing from  $< 1$  at time 0 to  $> 1$  at time  $t$ ;
2. HR is monotone increasing from  $< 1$  at time 0 to 1 at time  $t$ ;
3. HR is monotone decreasing from 1 at time 0 to  $< 1$  at time  $t$ .

The first setting corresponds to the worst violation of the PH assumption: crossing hazards. In this case, we let  $\alpha(u) = \log(0.46) + (u/t_M)^s$ , i.e., the HR is 0.46 at time zero, crosses 1 and eventually increases to 1.25 at time  $t_M$ . Here  $t_M$  is the longest follow-up time, which is greater than both  $t$  in  $t$ -MST and  $t_H$  in estimating the HR. In this case, the parameter  $s$  dictates how fast the HR increases with time and ARE. Although HR crosses one, two survival functions here don't cross within the interval  $[0, t_M]$ . The second setting corresponds to the case where the treatment benefit is large at time 0 but gradually diminishes. To characterize this pattern, we let  $\alpha(u) = -c(t_M - u)^s$ , where the constant  $c$  is chosen such that the average HR,  $t_M^{-1} \int_0^{t_M} \exp\{\alpha(u)\} du$ , is approximately 0.7. The last setting corresponds to the case, where the treatment benefit is small at the beginning but grows gradually during the follow-up. Specifically, we let  $\alpha(u) = -cu^s$ , where  $c$  is a constant such that the average HR is approximately 0.7. For all three settings, both the analytic ARE and ERE based on 5000 simulations are evaluated for all combinations of  $s \in \{0.5, 1, 2\}$ ,  $\lambda c \in \{0.06\%, 2.5\%\}$  and survival rates  $S_0(t) \in \{0.30, 0.50, 0.70\}$  and reported in Table 2. In the simulation, we let  $(t, t_H, t_M) = (40, 42.5, 43)$  to reflect the fact that HR often can be estimated over a larger time interval than  $t$ -MST.

When HR is monotone increasing and crosses one during the follow-up (Case 1), the  $t$ -MST based test can be substantially more efficient than the HR-based test; when the HR is less than one at time 0 and increases to 1 at the end the study (Case 2), the  $t$ -MST based test is also more efficient than the HR-based test but the difference is modest; and when the HR is 1 at time 0 and decreases over time (Case 3), the  $t$ -MST based test can be less efficient than the HR-based test. Similar to the case where the PH assumption holds, the EREs are in general consistent with their asymptotical counterparts in all settings, supporting the use of EREs for evaluating the finite sample performances of two tests.

## 5. Conclusions

In this paper, we compared the inference procedures based on HR and the  $t$ -MST difference under a hypothesis testing paradigm. Through an extensive numerical study with various treatment difference profiles, the AREs for the  $t$ -MST and HR tests are similar under the PH-models. For the non-PH models, unless the true difference of two survival functions is quite large near the end of the study and the censoring is light during the study follow-up, the  $t$ -MST test generally outperforms the HR test. The employed technique is analogous to those used for comparing the nonparametric tests for shift in location with the t-test when the normality assumption is violated (Ruberg, 1986). In summary, the  $t$ -MST based test is preferred if we are expecting an early benefit, for example, caused by less invasive treatment such as a low dose option. On the other hand, HR-based test is more powerful for detecting the delayed treatment effect encountered in recent oncology trials testing the efficacy of immunotherapy.

Perhaps a line of future research is to construct more powerful tests adaptive to the deviation from PH assumption. There are quite a few novel statistical tests for non-proportional hazards alternatives proposed and discussed in the literature (Fleming et al., 1987; Pepe and Fleming, 1989; Kosorok and Lin, 1999; Feng and Wahed, 2008; Yang and Prentice, 2010; Uno et al., 2015; Garès et al., 2015). However, there are no coherent corresponding estimation procedures to quantify the treatment difference. A more powerful test than the HR test for a non-PH case without a companion estimation procedure would have limited value clinically. To make coherent evaluation on the group difference, we suggest to first define a clinically meaningful measure summarizing the treatment benefit and employ the corresponding estimation and hypothesis testing procedure to conduct the statistical analysis. There are very few two-sample model-free estimating procedures for treatment effect, which can also be used for testing the equivalence of two survival functions in survival analysis. For example, one may use the difference or ratio of two  $t$ -year event rates, median survival times, and  $t$ -MSTs. Often the median event time is not estimable due to a short study followup time. Using the  $t$ -year event rate difference as the primary parameter of interest would ignore the temporal profile of the treatment effect before and after the time point  $t$ . The estimate for the  $t$ -MST difference or ratio appears to be a useful tool for analyzing event time observations for both its interpretability and the power of the associated hypothesis testing procedure.

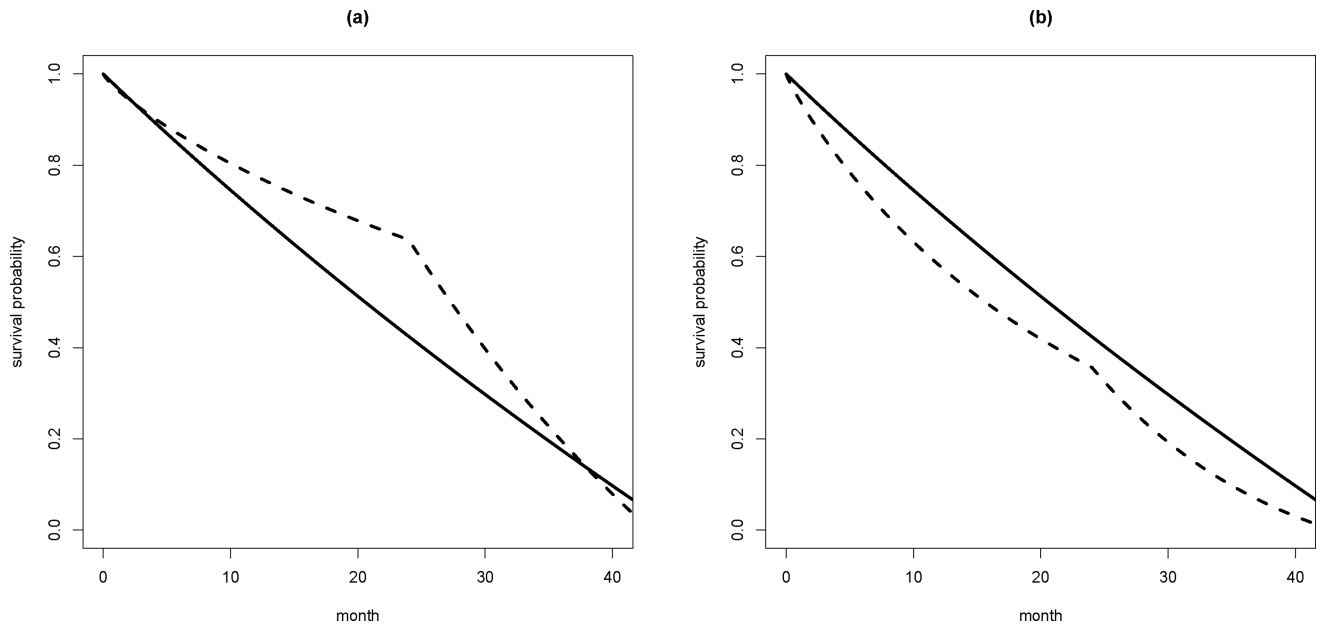
## Acknowledgments

This research was partially supported by R01 HL089778 (NIH/NHLBI), R00 HS022193 (NIH/AHRQ), and R21 AG049385 (NIH/NIA). We also would like to thank Editor, AE and a referee for their constructed comments to the paper.

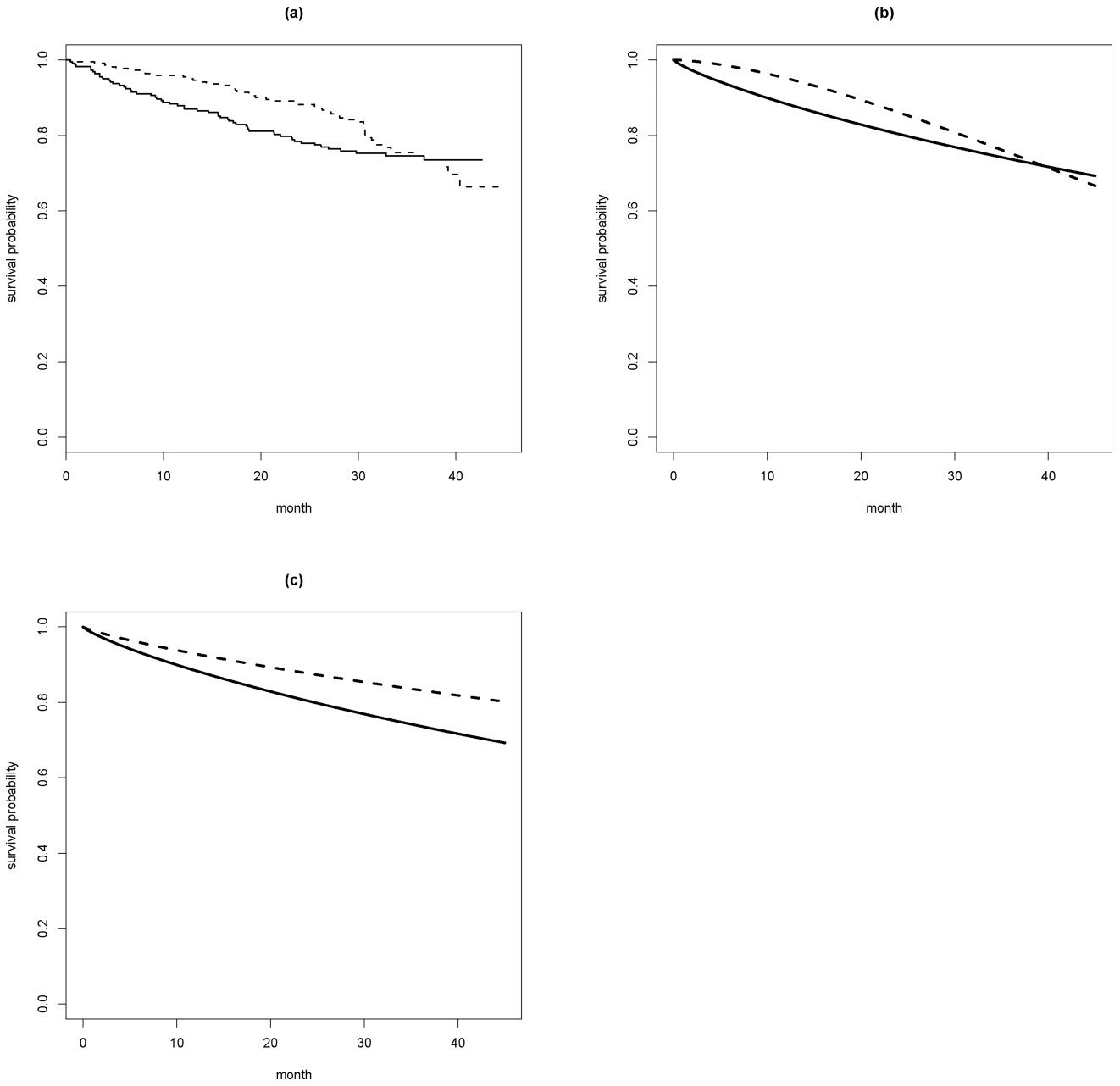
## References

- Cox D. Regression models and life-tables. *Journal of the Royal Statistical Society, Series B.* 1972; 34:187–220.
- Feng W, Wahed AS. Supremum weighted log-rank test and sample size for comparing two-stage adaptive treatment strategies. *Biometrika.* 2008; 95:695–707.
- Fleming, TR., Harrington, DP. *Counting Processes and Survival Analysis.* John Wiley & Sons; 2011.
- Fleming TR, Harrington DP, O'sullivan M. Supremum versions of the log-rank and generalized Wilcoxon statistics. *Journal of the American Statistical Association.* 1987; 82:312–320.
- Garès V, Andrieu S, Dupuy JF, Savy N. An omnibus test for several hazard alternatives in prevention randomized controlled clinical trials. *Statistics in Medicine.* 2015; 34:541–557. [PubMed: 25388274]
- Harrington D, Fleming T. A class of rank test procedures for censored survival data. *Biometrika.* 1982; 69:553–566.
- Hodges J, Lehmann E. The efficiency of some nonparametric competitors of the t-test. *Annals of Mathematical Statistics.* 1956; 27:324–335.
- Hoeffding W, Rosenblatt J. The efficiency of tests. *Annals of Mathematical Statistics.* 1955; 26:52–63.
- Kalbeisch J, Prentice R. Estimation of the average hazard ratio. *Biometrika.* 1981; 68:105–112.
- Karrison T. Restricted mean life with adjustment for covariates. *Journal of American Statistical Association.* 1987; 82:1169–1176.
- Kosorok MR, Lin CY. The versatility of function-indexed weighted log-rank statistics. *Journal of the American Statistical Association.* 1999; 94:320–332.
- Lagakos S. The loss in efficiency from misspecifying covariates in proportional hazards regression models. *Biometrika.* 1988; 75:156–160.

- Lin D, Wei LJ. The robust inference for the Cox proportional hazards model. *Journal of American Statistical Association*. 1989; 84:1074–1078.
- Pepe MS, Fleming TR. Weighted Kaplan-Meier statistics: a class of distance tests for censored survival data. *Biometrics*. 1989; 45:497–507. [PubMed: 2765634]
- Pepe MS, Fleming TR. Weighted Kaplan-Meier statistics: Large sample and optimality considerations. *Journal of the Royal Statistical Society, Series B*. 1991; 53:341–352.
- Rajkumar S, Jacobus S, Callander N, Fonseca R, Vesole D, Williams M, et al. Lenalidomide plus high-dose dexamethasone versus lenalidomide plus low-dose dexamethasone as initial therapy for newly diagnosed multiple myeloma: an open-label randomised controlled trial. *Lancet Oncology*. 2010; 11:29–37. [PubMed: 19853510]
- Royston P, Parmar M. The use of restricted mean survival time to estimate the treatment effect in randomized clinical trials when the proportional hazards assumption is in doubt. *Journal of American Statistical Association*. 2011; 30:2409–2421.
- Ruberg SJ. Efficiencies of some two-sample location tests for a broad class of distributions. *Communications in Statistics - Theory and Methods*. 1986; 15:2991–3004.
- Schoenfeld D. The asymptotic properties of nonparametric tests for comparing survival distributions. *Biometrika*. 1981; 68:316–319.
- Schoenfeld D, Richter J. Nomograms for calculating the number of patients needed for a clinical trial with survival as an endpoint. *Biometrics*. 1982; 38:163–170. [PubMed: 7082758]
- Slud E. Relative efficiency of the logrank test within a multiplicative intensity model. *Biometrika*. 1991; 18:1172–1187.
- Tian L, Zhao L, Wei LJ. Predicting the restricted mean event time with the subject's baseline covariates in survival analysis. *Biostatistics*. 2014; 15:222–233. [PubMed: 24292992]
- Trinquart L, Jacot J, Conner SC, Porcher R. Comparison of treatment effects measured by the hazard ratio and by the ratio of restricted mean survival times in oncology randomized controlled trials. *Journal of Clinical Oncology*. 2016; 34:1813–1819. [PubMed: 26884584]
- Uno H, Claggett B, Tian L, Inoue E, Gallo P, Miyata T, et al. Moving beyond the hazard ratio in quantifying the between-group difference in survival analysis. *Journal of Clinical Oncology*. 2014; 32:2380–2385. [PubMed: 24982461]
- Uno H, Tian L, Claggett B, Wei LJ. A versatile test for equality of two survival functions based on weighted differences of Kaplan–Meier curves. *Statistics in Medicine*. 2015; 34:3680–3695. [PubMed: 26194988]
- Wei LJ. The accelerated failure time model: A useful alternative to the Cox regression model in survival analysis. *Statistics in Medicine*. 1992; 11:1871–1875. [PubMed: 1480879]
- Yang S, Prentice R. Improved logrank-type tests for survival data using adaptive weights. *Biometrics*. 2010; 66:30–38. [PubMed: 19397582]
- Zhang D, Quan H. Power and sample size calculation for log-rank test with a time lag in treatment effect. *Statistics in Medicine*. 2009; 28:864–879. [PubMed: 19152230]
- Zhao L, Tian L, Uno H, Solomon S, Pfeffer M, Schindler J, et al. Utilizing the integrated difference of two survival functions to quantify the treatment contrast for designing, monitoring, and analyzing a comparative clinical study. *Clinical Trials*. 2012; 9:570–577. [PubMed: 22914867]
- Zucker D. Restricted mean life with covariates: modification and extension of a useful survival analysis method. *Journal of American Statistical Association*. 1998; 93:702–709.



**Figure 1.** The weight functions for  $t$ -MST and HR based tests; solid: weight function for  $t$ -MST based test; dotted: weight function for HR based test; (a) light drop-out rate with  $\lambda_c = 0.06\%$  and (b) high drop-out rate with  $\lambda_c = 2.5\%$



**Figure 2.** (a) The KM curves in two arms of the ECOG-ACRIN study; (b) The survival curves based on the Weibull model in two arms of the ECOG-ACRIN study without assuming the PH assumption; (c) The survival curves based on the Weibull model in the high dose arm of the ECOG-ACRIN study assuming the PH assumption with a HR of 0.7; solid: high dose arm; dotted: low dose arm.

**Table 1**

Asymptotical relative efficiency (ARE) and empirical relative efficiency (ERE) under PH alternatives with a HR of 0.7; EREs are estimated based on 5000 sets simulated data.

Censoring	ARE(ERE)	
	light	heavy
$S_0(t) = 0.90$	0.95(0.94)	1.05(1.03)
$S_0(t) = 0.80$	0.96(0.96)	1.05(1.03)
$S_0(t) = 0.70$	0.97(0.93)	1.05(1.01)
$S_0(t) = 0.60$	0.98(0.96)	1.06(1.03)
$S_0(t) = 0.50$	0.99(0.96)	1.06(1.02)
$S_0(t) = 0.40$	1.01(0.97)	1.06(1.02)
$S_0(t) = 0.30$	1.02(0.97)	1.06(1.02)
$S_0(t) = 0.20$	1.04(0.99)	1.05(1.02)
$S_0(t) = 0.10$	1.05(1.01)	1.04(1.01)

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript



Asymptotical relative efficiency (ARE) and empirical relative efficiency (ERE) under nonproportional hazards alternatives; EREs are estimated based on 5000 sets simulated data.

**Table 2**

Censoring	$S_0(t) = 0.30$			$S_0(t) = 0.50$			$S_0(t) = 0.70$		
	ARE(ERE)		heavy	ARE(ERE)		heavy	ARE(ERE)		heavy
	light	heavy		light	heavy		light	heavy	
	$\alpha(w) = \log(0.46) + (w/t)^s$								
$s = 0.5$	1.76(1.77)	1.16(1.13)	1.16(1.13)	2.02(2.05)	1.24(1.22)	1.24(1.22)	2.30(2.37)	1.34(1.31)	1.34(1.31)
$s = 1.0$	1.34(1.33)	1.11(1.08)	1.11(1.08)	1.42(1.41)	1.15(1.13)	1.15(1.13)	1.49(1.47)	1.20(1.17)	1.20(1.17)
$s = 2.0$	1.18(1.17)	1.09(1.07)	1.09(1.07)	1.21(1.19)	1.11(1.10)	1.11(1.10)	1.23(1.23)	1.13(1.12)	1.13(1.12)
	$\alpha(w) = -\alpha(\alpha)(t-w)^s$								
$s = 0.5$	1.10(1.05)	1.07(1.05)	1.07(1.05)	1.10(1.06)	1.09(1.06)	1.09(1.06)	1.10(1.07)	1.10(1.06)	1.10(1.06)
$s = 1.0$	1.20(1.17)	1.09(1.06)	1.09(1.06)	1.22(1.19)	1.11(1.09)	1.11(1.09)	1.24(1.21)	1.13(1.11)	1.13(1.11)
$s = 2.0$	1.37(1.34)	1.11(1.09)	1.11(1.09)	1.43(1.38)	1.16(1.13)	1.16(1.13)	1.48(1.45)	1.19(1.17)	1.19(1.17)
	$\alpha(w) = -\alpha(\alpha)w^s$								
$s = 0.5$	0.84(0.80)	1.02(0.99)	1.02(0.99)	0.81(0.77)	1.00(0.96)	1.00(0.96)	0.78(0.76)	0.98(0.94)	0.98(0.94)
$s = 1.0$	0.74(0.70)	1.00(0.96)	1.00(0.96)	0.70(0.67)	0.96(0.92)	0.96(0.92)	0.67(0.64)	0.92(0.90)	0.92(0.90)
$s = 2.0$	0.61(0.58)	0.96(0.93)	0.96(0.93)	0.56(0.54)	0.89(0.86)	0.89(0.86)	0.53(0.52)	0.84(0.82)	0.84(0.82)