

2024 Summer Institute In Statistics for Clinical & Epidemiological Research

Module 3:

Design, Conduct, and Analysis of Randomized Clinical Trials with Time to Event Primary Endpoints

.....

Lecture 28:

Adaptive RCT Designs with Time to Event Endpoints

Scott S. Emerson, M.D., Ph.D.
Professor Emeritus of Biostatistics
University of Washington

1

Where Am I Going?

.....

Overview and Organization of the Talk

2

2

Science and Statistics



- Statistics is about science
 - (Science in the broadest sense of the word)
- Science is about proving things to people
 - (The validity of any proof rests solely on the willingness of the audience to believe it)
- In RCT, we are trying to prove the effect of some treatment
 - What do we need to consider as we strive to meet the burden of proof with adaptive modification of a RCT design?
- Does time to event data affect those issues?
 - Short answer: No, UNLESS subject to censoring
 - So, true answer: Yes.

3

3

Overview: Sequential, Adaptive RCT



- Increasing interest in the use of sequential, adaptive RCT designs
 - More efficient “drug discovery” for “personalized medicine”
 - More ethical treatment of individuals and populations
- FDA Draft guidance on adaptive designs
 - Well understood methods
 - Fixed sample
 - Group sequential
 - Blinded adaptation
 - Less well understood methods
 - Adaptive sample size re-estimation
 - Adaptive enrichment
 - Response-adaptive randomization
 - Adaptive selection of doses and/or treatments

4

4

Overview: Time-to-Event



- Many confirmatory phase 3 RCTs compare the distribution of time to some event (e.g., time to death or progression free survival).
- Common statistical analyses: Logrank test and/or PH regression
- Just as commonly: True distributions do not satisfy PH
- Providing users are aware of the nuances of those methods, such departures need not preclude the use of those methods

5

5

Overview: Premise



- Much of the concern with “less well understood” methods has to do with “less well understood” aspects of survival analysis in RCT
 - *“Everyone is ignorant, just on different subjects” – Will Rogers*
- Proportional hazards holds under strong null
 - But weak null can be important (e.g., noninferiority)
- Log linear hazard may be close to linear in log time over support of censoring distribution → approximately Weibull
 - A special case of PH only when shape parameter is constant
- Hazard ratio estimate can be thought of a weighted time-average of ratio of hazard functions
 - But in Cox regression, weights depend on censoring distribution
 - And in sequential RCT, censoring distribution keeps changing

6

6

Organization of the Presentation



- Sequential methods: Design and inference
 - Fixed sample designs
 - Group sequential designs
 - Adaptive designs
- **Adaptive methods with analyses of times to events**
 - **Efficiency of adaptive designs**
 - **Sample size re-estimation (SSRE) with low event rates and/or extreme effects**
 - Adaptive designs in presence of time-varying treatment effects
 - Potential impact of surrogate data

7

7

Adaptive RCT: Issues for Another Day



- Adaptive randomization ratios
 - Avoiding the introduction of confounding
- Adaptive enrichment
 - Designs and inference
- Operational issues

8

8

Abridged Version

.....

“He was against it.”

- Calvin Coolidge

9

9

Sequential Methods

.....

[How do “Adaptive Designs” differ from previously described methods?](#)

Where am I going?
I present some examples where the behavior of standard analysis methods for time-to-event data are not well understood

10

10

RCT Phases of Investigation

- A sequential, adaptive process
 - But only “piecewise continuous”
- During any individual clinical trial
 - Sequential monitoring, adaptation addresses that trial’s issues
- “White space” between trials: Detailed and exploratory analyses
 - Evaluation of multiple endpoints; cost/benefit tradeoffs
 - Exploratory analyses
 - Integration of results from other studies
 - Management decisions
 - Regulatory and ethical review
- Next RCT (if any): May address different question or indication

11

11

Science: Treatment “Indication”

- Disease
 - Therapy: Putative cause vs signs / symptoms
 - May involve method of diagnosis, response to therapies
 - Prevention / Diagnosis: Risk classification
- Population
 - Therapy: Restrict by risk of AEs or actual prior experience
 - Prevention / Diagnosis: Restrict by contraindications
- Treatment or treatment strategy
 - Formulation, administration, dose, frequency, duration, ancillary therapies
- Outcome
 - Clinical vs surrogate; timeframe; method of measurement

12

12

Notation

.....

Baseline data :	$W_1, W_2, W_3, \dots, W_N$
Treatment data :	$X_1, X_2, X_3, \dots, X_N$
Potential data :	$Y_1, Y_2, Y_3, \dots, Y_N$
Probability model :	$Y_i X_i, W_i \stackrel{ind}{\sim} F_i$
Target of inference :	$\theta = \theta(F_1, \dots, F_N)$
Estimated treatment effect :	$\hat{\theta}_N = \theta(\hat{F}_1, \dots, \hat{F}_N) \sim N(\theta, V(\theta)/N)$
Normalized test statistic :	$Z_N = \frac{\hat{\theta}_N - \theta_0}{\sqrt{V(\theta_0)/N}} \sim N\left(\frac{\theta - \theta_0}{\sqrt{V(\theta_0)/N}}, 1\right)$
P value :	$P_N = \Phi(Z_N) \stackrel{H_0: \theta = \theta_0}{\sim} U(0, 1)$

13

13

Inference

.....

- At the end of the study, report four numbers

- Frequentist
 - Estimated treatment effect (low bias, consistent, low MSE)
 - Confidence interval (a counterfactual: hypotheses leading to data)
 - P value (perhaps a counterfactual)

- Bayesian (for some prior or population of priors)
 - Summary measure for posterior distribution
 - Credible interval
 - Posterior probability of relevant hypotheses

14

14

Clinical Trial Design

- Design the study to discriminate between important hypotheses
 - Confidence / credible intervals should not contain both of two competing hypotheses
- Finding an approach that best addresses the often competing goals: Science, Ethics, Efficiency
 - Basic scientists: focus on mechanisms
 - Clinical scientists: focus on overall patient health
 - Ethical: focus on patients on trial, future patients
 - Economic: focus on profits and/or costs
 - Governmental: focus on safety of public: treatment safety, efficacy, marketing claims
 - Statistical: focus on questions answered precisely
 - Operational: focus on feasibility of mounting trial

15

15

Design: Distinctions without Differences

- There is no such thing as a “Bayesian design”
- Every RCT design has a Bayesian interpretation
 - (And each person may have a different such interpretation)
- Every RCT design has a frequentist interpretation
 - (In poorly designed trials, this may not be known exactly)
- I focus on the use of both interpretations
 - Phase 2: Bayesian probability space
 - Phase 3: Frequentist probability space
 - Entire process: Both Bayesian and frequentist optimality criteria

16

16

Application to Drug Discovery

- We consider a population of candidate drugs
- We use RCT to “diagnose” truly beneficial drugs
- Use both frequentist and Bayesian optimality criteria
 - Sponsor:
 - High probability of adopting a beneficial drug (frequentist power)
 - Regulatory:
 - Low probability of adopting ineffective drug (freq type 1 error)
 - High probability that adopted drugs work (posterior probability)
 - Public Health (frequentist sample space, Bayes criteria)
 - Maximize the number of good drugs adopted
 - Minimize the number of ineffective drugs adopted

17

17

Frequentist vs Bayesian: Bayes Factor

- Frequentist and Bayesian inference truly complementary
 - Frequentist: Design so the same data not likely from null / alt
 - Bayesian: Explore updated beliefs based on a range of priors
- Bayes rule tells us that we can parameterize the positive predictive value by the type I error and prevalence
 - Maximize new information by maximizing Bayes factor
 - With simple hypotheses:

$$PPV = \frac{\text{power} \times \text{prevalence}}{\text{power} \times \text{prevalence} + \text{type I err} \times (1 - \text{prevalence})}$$

$$\frac{PPV}{1 - PPV} = \frac{\text{power}}{\text{type I err}} \times \frac{\text{prevalence}}{1 - \text{prevalence}}$$

$$\text{posterior odds} = \text{Bayes Factor} \times \text{prior odds}$$

18

18

Sequential RCT



- Ethical and efficiency concerns can be addressed through sequential sampling
- During the conduct of the study, data are analyzed at periodic intervals and reviewed by the DMC
- Using interim estimates of treatment effect decide whether to continue the trial
- If continuing, decide on any modifications to
 - scientific / statistical hypotheses and/or
 - sampling scheme

19

19

Notation: Sampling Independent Groups



Independent groups :	$j = 1, \dots, J$
Baseline data :	$W_{j1}, \dots, W_{j\tilde{N}_j}$
Treatment data :	$X_{j1}, \dots, X_{j\tilde{N}_j}$
Potential data :	$Y_{j1}, \dots, Y_{j\tilde{N}_j}$
Probability model :	$Y_{ji} X_{ji}, W_{ji} \stackrel{ind}{\sim} F_{ji}$
Target of inference :	$\tilde{\theta}_j = \tilde{\theta}_j(F_{j1}, \dots, F_{j\tilde{N}_j})$
Estimated treatment effect :	$\hat{\tilde{\theta}}_{j\tilde{N}_j} = \tilde{\theta}_j(\hat{F}_{j1}, \dots, \hat{F}_{j\tilde{N}_j}) \sim N(\tilde{\theta}_j, V_j(\tilde{\theta}_j)/\tilde{N}_j)$
Normalized test statistic :	$\tilde{Z}_{j\tilde{N}_j} = \frac{\hat{\tilde{\theta}}_{j\tilde{N}_j} - \tilde{\theta}_{j0}}{\sqrt{V_j(\tilde{\theta}_{j0})/\tilde{N}_j}} \sim N\left(\frac{\tilde{\theta}_j - \tilde{\theta}_{j0}}{\sqrt{V_j(\tilde{\theta}_{j0})/\tilde{N}_j}}, 1\right)$
P value :	$\tilde{P}_{j\tilde{N}_j} = \Phi(\tilde{Z}_{j\tilde{N}_j}) \stackrel{H_{j0}: \tilde{\theta}_j = \tilde{\theta}_{j0}}{\sim} U(0, 1)$

20

20

Adaptive RCT Design: Precursors



- What
 - Sequential Probability Ratio Test (Wald, classified during WWII)
 - Group sequential designs (Armitage, et al., 1969)
 - Bayesian methods
- Why
 - [Sample size adjustment](#)
 - Selection of treatment arms
 - Selection of subgroups
 - Multiple endpoints
- How
 - Frequentist sampling plans
 - Prespecified maximum statistical information, sampling plan
 - Control of frequentist experimentwise type 1 error, power
 - Sequential inference generally well-described
 - Selection of Bayesian priors

21

21

Notation: Group Sequential Designs



- A common treatment effect across groups
- Group size independent of prior estimates of treatment effect

Prespecified (rule for)	N_1, N_2, \dots, N_J
Potential data :	$Y_1, Y_2, Y_3, \dots, Y_{N_j}$
Probability model :	$Y_i \stackrel{iid}{\sim} (\theta, V)$
Interim estimates :	$\hat{\theta}_{N_j} = \hat{\theta}(Y_1, \dots, Y_{N_j})$
Without sequential sampling :	
Approximate distn :	$\hat{\theta}_j = \hat{\theta}_{N_j} \sim N(\theta, V / N_j)$
Indep increments :	$Cov(\hat{\theta}_{N_j}, \hat{\theta}_{N_{j+1}}) = V / N_{j+1}$
Interim test statistics :	$Z_j = Z_{N_j} = \frac{\hat{\theta}_j - \theta_0}{\sqrt{V / N_j}}$

22

22

Group Sequential Designs

- Perform analyses when sample sizes N_1, \dots, N_j
 - Can be randomly determined
- At each analysis choose stopping boundaries
 - $a_j < b_j < c_j < d_j$
- Compute test statistic $T_j = T(X_1, \dots, X_{N_j})$
 - Stop if $T_j < a_j$ (extremely low)
 - Stop if $b_j < T_j < c_j$ (approximate equivalence)
 - Stop if $T_j > d_j$ (extremely high)
 - Otherwise continue
- Boundaries chosen to protect 2 of 3 operating characteristics
 - Type 1 error, power
 - Type 1 error, power, maximal sample size

23

23

Distinctions without Differences

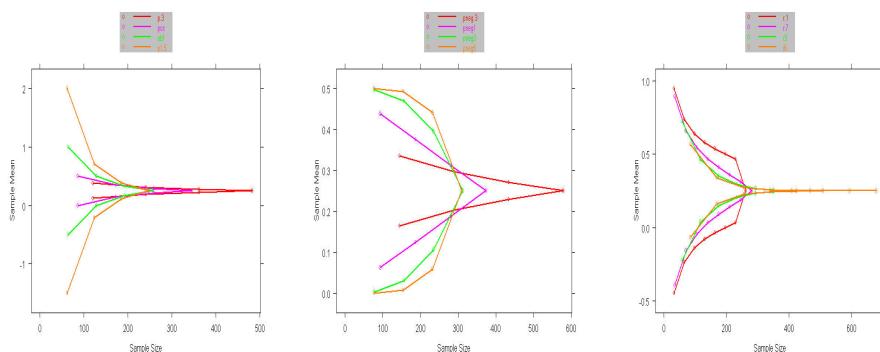
- Boundary scales (1:1 transformations among these)
 - Z statistic
 - P value
 - Fixed sample (so wrong)
 - Computed under sequential sampling rule (so correct)
 - Error spending function
 - Estimates
 - MLE (biased due to stopping rule)
 - Adjusted for stopping rule
 - Conditional power
 - Computed under design alternative
 - Computed under current MLE
 - Predictive power
 - Computed under flat prior (possibly improper)
 - Bayesian posterior probabilities

24

24

Spectrum of Boundary Shapes

- All of the rules depicted have the same type I error and power to detect the design alternative



25

Operating Characteristics

- For any pre-specified stopping rule, however, we can compute the correct sampling distribution with specialized software
- From the computed sampling distributions we then compute
 - Bias adjusted estimates
 - Correct (adjusted) confidence intervals
 - Correct (adjusted) P values
- Candidate designs are then compared with respect to their operating characteristics

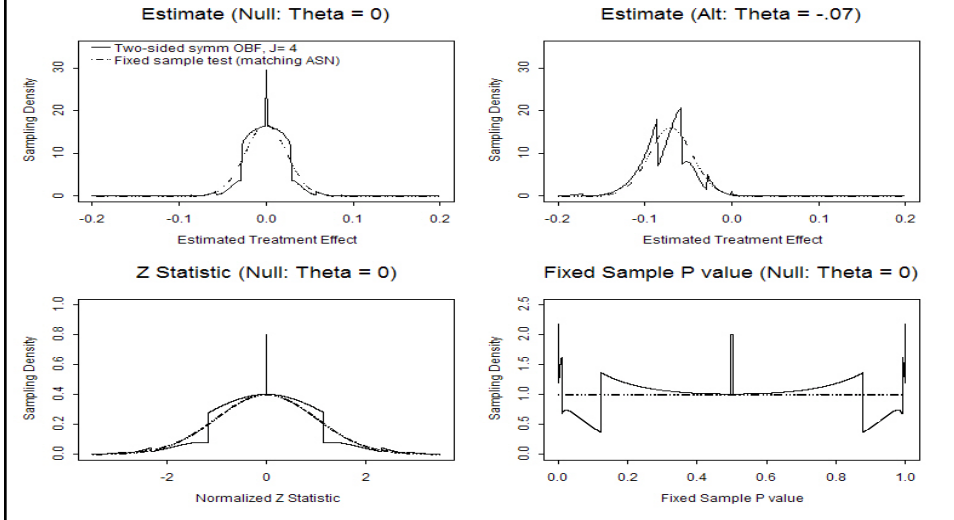
26

26

Sampling Densities: Estimate, Z, Fixed P



- For a particular stopping rule



27

Evaluation of Designs



- Process of choosing a trial design
 - Define candidate design
 - Usually constrain two operating characteristics
 - Type I error, power at design alternative
 - Type I error, maximal sample size
 - Evaluate other operating characteristics
 - Different criteria of interest to different investigators
 - Modify design
 - Iterate

28

28

Which Operating Characteristics



- The same regardless of the type of stopping rule
- Frequentist power curve
 - Type I error (null) and power (design alternative)
- Sample size requirements
 - Maximum, average, median, other quantiles
 - Stopping probabilities
 - Tradeoffs between sample size and power
- Inference at study termination (at each boundary)
 - Frequentist or Bayesian (under spectrum of priors)
- (Futility measures
 - Conditional power, predictive power)

29

29

But What If ...?



- Possible motivations for adaptive designs
- Changing conditions in medical environment
 - Approval / withdrawal of competing / ancillary treatments
 - Diagnostic procedures
- New knowledge from other trials about similar treatments
- Evidence from ongoing trial
 - Toxicity profile (therapeutic index)
 - Interim estimates of primary efficacy / effectiveness endpoint
 - Overall
 - Within subgroups
 - Interim alternative analyses of primary endpoints
 - Interim estimates of secondary efficacy / effectiveness endpoints

30

30

“Modern” Adaptive RCT Design



- What if maximal statistical information and sampling plan is not prespecified?
- First: Control of experimentwise type 1 error
 - (Generate a random uniform in a dark room)
 - Use Fisher’s combination of p values (Bauer & Koehne, 1994)
 - Conditional error functions (Proschan & Hunsberger, 1995)
 - Re-weighted incremental statistics (Fisher, 1998; Cui, et al, 1999)
 - “Bayesian adaptive designs” based on predictive probabilities with simulations to verify control of type 1 error

31

31

Adaptive Sampling: General Case



- At each interim analysis, possibly modify statistical or scientific aspects of the RCT
- Primarily statistical characteristics
 - Maximal statistical information (UNLESS: impact on MCID)
 - Schedule of analyses (UNLESS: time-varying effects)
 - Conditions for stopping (UNLESS: time-varying effects)
 - Randomization ratios (UNLESS: introduce confounding)
 - Statistical criteria for credible evidence
- Primarily scientific characteristics
 - Target patient population (inclusion, exclusion criteria)
 - Treatment (dose, administration, frequency, duration)
 - Clinical outcome and/or statistical summary measure

32

32

Adaptive Sampling: Issues

- How do the newer adaptive approaches relate to the constraint of human experimentation and scientific method?
- Effect of adaptive sampling on trial ethics and efficiency
 - Avoiding unnecessarily exposing subjects to inferior treatments
 - Avoiding unnecessarily inflating the costs (time / money) of RCT
- Effect of adaptive sampling on scientific interpretation
 - Exploratory vs confirmatory clinical trials
- Effect of adaptive sampling on statistical credibility
 - Control of type I error in frequentist analyses
 - Promoting predictive value of “positive” trial results

33

33

Typical Adaptive Design

- Perform analyses when sample sizes N_1, \dots, N_J
 - Can be randomly determined
- At each analysis choose stopping boundaries
 - $a_j < b_j < c_j < d_j$
- Compute test statistic $T_j = T(X_1, \dots, X_{N_j})$
 - Stop if $T_j < a_j$ (extremely low)
 - Stop if $b_j < T_j < c_j$ (approximate equivalence)
 - Stop if $T_j > d_j$ (extremely high)
 - Otherwise continue
- At penultimate analysis ($J-1$), use unblinded interim test statistic to choose final sample size N_J or to modify other aspects of RCT

34

34

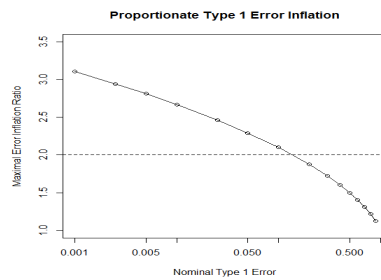
Proschan & Hunsberger



- Worst case type I error of two stage design

$$\alpha_{worst} = 1 - \Phi(a_2^{(Z)}) + \frac{\exp\left(-\left(a_2^{(Z)}\right)^2 / 2\right)}{4},$$

- Can be more than two times the nominal
 - $a_2 = 1.96$ gives type I error of 0.0616
 - (Compare to Bonferroni results)



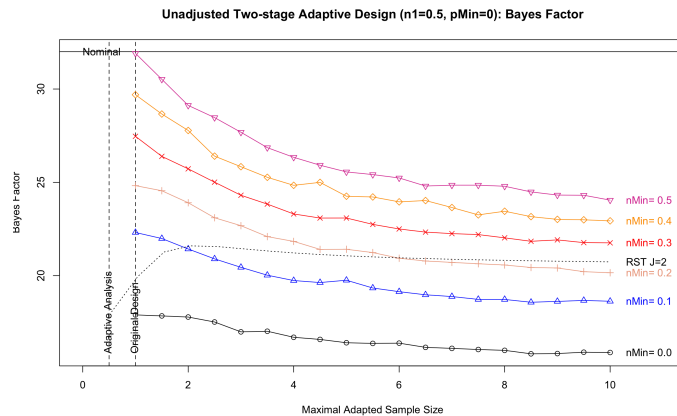
35

35

Modified Adaptive Rules: Bayes Factor



- Second stage sample size
 - Minimal sample size if $Z_1 < 0$ (worst case used infinity)
 - Minimal sample size if $Z_1 > z_{1-\alpha/2}$ (worst case used 0)
 - Bounded maximal sample size (many use 2-fold increase)



36

36

Adaptive Control of Type 1 Errors

- Proschan and Hunsberger (1995)
 - Adaptive modification of RCT design at a single interim analysis can more than double type 1 error unless carefully controlled
- Those authors describe adaptations to maintain experimentwise type I error and increase conditional power
 - Must prespecify a conditional error function

$$\int_{-\infty}^{\infty} A(z) \phi(z) dz = \alpha.$$

- Often choose function from some specified test

$$A(z) = Pr_{\delta=0}(Z_2 \geq \Phi^{-1}(1 - \alpha) | \tilde{Z}_1 = z, \tilde{n}_2 = n_2 - n_1),$$

- Find critical value to maintain type I error

$$Pr_{\delta=0}(Z_2^* \geq c(\tilde{n}_2^*, \tilde{z}_1) | \tilde{n}_2^*(\tilde{z}_1)) = A(\tilde{z}_1).$$

37

37

Conditional Distn: Immediate Outcomes

- Sample size N_j and parameter θ_j can be adaptively chosen based on data from prior stages $1, \dots, j-1$
 - (Most often we choose $\theta_j = \theta$ with immediate data)

$$\hat{\theta}_j | \tilde{N}_j \sim N\left(\theta, \frac{V(\theta)}{\tilde{N}_j}\right)$$

$$\tilde{Z}_j | \tilde{N}_j \sim N\left(\frac{\hat{\theta}_j - \theta_0}{\sqrt{V(\theta)/\tilde{N}_j}}, 1\right)$$

$$\tilde{P}_j | \tilde{N}_j \stackrel{H_0}{\sim} U(0, 1).$$

Conditional distributions
 are totally independent
 under the null hypothesis

38

38

Without Adaptation

.....

- Statistic at the j -th analysis a weighted average of data accrued between analyses

Statistics computed on j th increment : $\hat{\theta}_j$ \tilde{Z}_j \tilde{P}_j

$$\hat{\theta}_j = \frac{\sum_{k=1}^j \tilde{N}_k \hat{\theta}_k}{N_j} \qquad Z_j = \frac{\sum_{k=1}^j \sqrt{\tilde{N}_k} \tilde{Z}_k}{\sqrt{N_j}}.$$

39

39

Protecting Type I Error

.....

- LD Fisher's variance spending method
 - Arbitrary hypotheses $H_{0j}; \theta_j = \theta_{0j}$
 - Incremental test statistics Z_j^*
 - Allow arbitrary weights W_j specified at stage $j-1$

$$Z_J = \frac{\sum_{k=1}^J \sqrt{W_k} \tilde{Z}_k}{\sqrt{\sum_{k=1}^J W_k}}$$

- RA Fisher's combination of P_j values (Bauer & Köhne)

$$P_J = \prod_{k=1}^J \tilde{P}_k$$

40

40

Unconditional Distribution



- Under the null
 - SDCT: Standard normal
 - Bauer & Kohne: Sum of exponentials
- Under the alternative
 - Unknown unless prespecified adaptations

$$\Pr(\tilde{Z}_j \leq z) = \sum_{n=0}^{\infty} \Pr(\tilde{Z}_j \leq z \mid \tilde{N}_j) \Pr(\tilde{N}_j = n).$$

41

41

Approaches for Testing



- If modify sample size at second stage (Cui, Hung, & Wang)

$$\tilde{N}_2^* = \tilde{N}_2^*(\tilde{Z}_1) \quad \tilde{Z}_2^* \text{ incremental statistic with revised } \tilde{N}_2^*$$

$$Z_2^* = \sqrt{\frac{\tilde{N}_1}{N_2}} \tilde{Z}_1 + \sqrt{\frac{\tilde{N}_2^*}{N_2}} \tilde{Z}_2^* \stackrel{H_0}{\sim} N(0,1)$$

- Equivalently, calculate Z statistic as usual and use different critical value

$$\text{reject } H_0 \Leftrightarrow Z_2^* = \sqrt{\frac{\tilde{N}_1}{N_2}} \tilde{Z}_1 + \sqrt{\frac{\tilde{N}_2^*}{N_2}} \tilde{Z}_2^* > b(\tilde{Z}_1, \tilde{N}_2^*)$$

$$b(\tilde{Z}_1, \tilde{N}_2^*) = \frac{1}{\sqrt{\tilde{N}_2^*}} \left[\sqrt{\frac{\tilde{N}_2^*}{N_2}} (z_{1-\alpha} \sqrt{N_2} - Z_1 \sqrt{\tilde{N}_1}) + Z_1 \sqrt{\tilde{N}_1} \right]$$

42

42

Choice of Adaptive Rule



Sample Size Re-estimation (SSRE)

Where am I going?

Some investigators desire to modify sample size more flexibly than allowed with GST

43

43

Comments



- In order to use these methods, we must have
 - The Z statistic at the adaptive analysis
 - The information growth function
 - Some adaptive rule (pre-specified or unspecified)
- It is easily shown that a minimal sufficient statistic is (Z, N) at stopping
 - All methods advocated for type 1 error control with fully adaptive designs are thus not based on sufficient statistics
 - Instead they re-weight data after the adaptive analysis
 - Changing the critical value is equivalent to re-weighted data!
- If the adaptive rule is not pre-specified, we must protect ourselves against everything

44

44

Fully Adaptive Sampling Plans



“Keep an open mind, but not so open that your brains fall out.”

- Virginia Gildersleeve?

45

45

Example



Adaptive Increase in Sample Size when Interim Results are Promising: A Practical Guide with Examples

Cyrus R. Mehta^{1,2}, Stuart J. Pocock³

¹Cytel Corporation, ²Harvard School of Public Health, ³London School of Hygiene and Tropical Medicine

SUMMARY

This paper discusses the benefits and limitations of adaptive sample size re-estimation for phase 3 confirmatory clinical trials. Comparisons are made with more traditional fixed sample and group sequential designs. It is seen that the real benefit of the adaptive approach arises through the ability to invest sample size resources into the trial in stages. The trial starts with a small up-front sample size commitment. Additional sample size resources are committed to the trial only if promising results are obtained at an interim analysis. This strategy is shown through examples of actual trials, one in neurology and one in cardiology, to be more advantageous than the fixed sample or group sequential approaches in certain settings. A major factor that has generated controversy and inhibited more widespread use of these methods has been their reliance on non-standard tests and p-values for preserving the type-1 error. If, however, the sample size is only increased when interim results are promising, one can dispense with these non-standard methods of inference. Therefore, in the spirit of making adaptive increases in trial size more widely appealing and readily implementable we here define those promising circumstances in which a conventional final inference can be performed while preserving the overall type-1 error. Methodological, regulatory and operational issues are examined. Copyright © 2000 John Wiley & Sons, Ltd.

46

46

(Counter) Example

.....

Commentary

Received 3 March 2011, Accepted 23 March 2011, Published online in Wiley Online Library
 (wileyonlinelibrary.com) DOI: 10.1002/sim.4271

**Statistics
in Medicine**

Comments on ‘Adaptive increase in sample size when interim results are promising: A practical guide with examples’

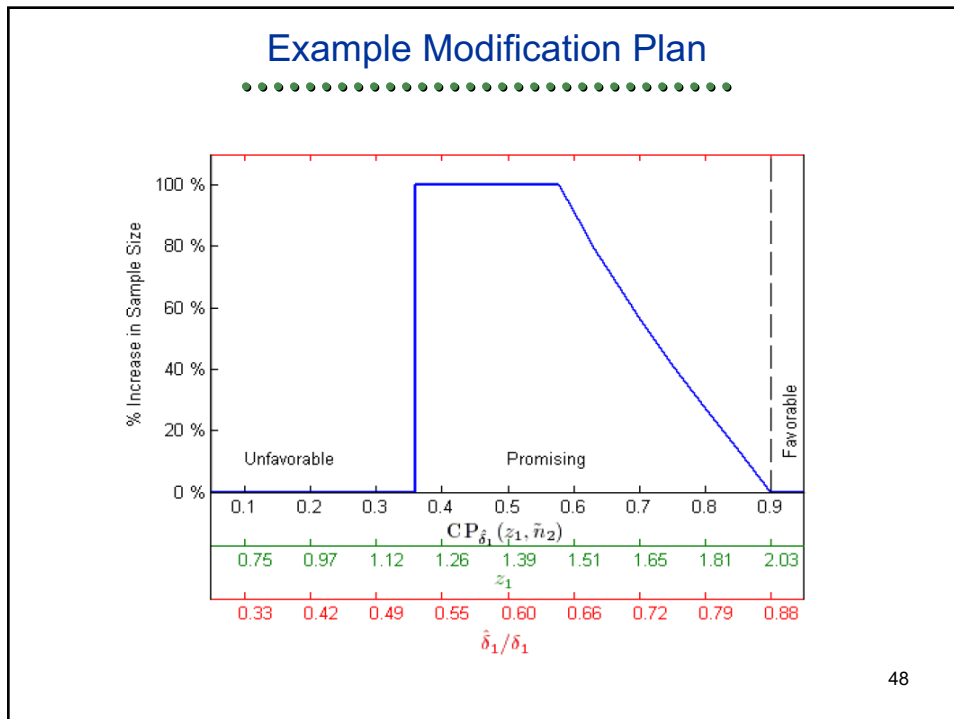
Scott S. Emerson,^{a*†} Gregory P. Levin^a and Sarah C. Emerson^b

Keywords: adaptive design; clinical trial; group sequential test; group sequential trial; statistical efficiency

In their paper [1], Drs Mehta and Pocock illustrate the use of a particular approach to revising the maximal sample size of a randomized clinical trial (RCT) by using an interim estimate of the treatment effect. Slightly extending the results of Gao *et al.* [2], the authors define conditions on an adaptive rule such that one can know that the naive statistical hypothesis test that ignores the adaptation is conservative. They then use this knowledge to define an adaptive rule for a clinical trial. In our review of this paper, however, we do not find that such an adaptive rule confers any advantage by the usual criteria for clinical trial design. Rather, we find that the designs proposed in this paper are markedly inferior to alternative designs that the authors do not (but should) consider.

47

47



48

48

Comparisons Unconditional Power

Table IV. Operating Characteristics of Fixed Sample and Adaptive Designs

Value of δ	Fixed Sample Design		Plan 4 (Adaptive)	
	Power	Expected SampleSize	Power	Expected Sample Size
1.6	61%	442	65%	499
1.7	66%	442	71%	498
1.8	71%	442	75%	497
1.9	76%	442	79%	494
2.0	80%	442	83%	491
All Plan 4 results are based on 100,000 simulated trials				

49

49

Comparisons Conditional Power

Table V. Operating Characteristics of the Fixed Sample and Adaptive Designs, Conditional on Interim Outcome

δ	Interim Outcome	Probability of Interim Outcome	Power Conditional on Interim Outcome		Expected Sample Size	
			Fixed	Adaptive	Fixed	Adaptive
1.6	Unfavorable	36%	30%	30%	442	442
	Promising	23%	62%	82%	442	687
	Favorable	41%	87%	87%	442	442
1.7	Unfavorable	32%	34%	34%	442	442
	Promising	23%	67%	85%	442	685
	Favorable	45%	89%	89%	442	442
1.8	Unfavorable	29%	38%	38%	442	442
	Promising	23%	70%	88%	442	682
	Favorable	49%	91%	91%	442	442
1.9	Unfavorable	26%	43%	43%	442	442
	Promising	22%	74%	90%	442	679
	Favorable	52%	93%	93%	442	442
2.0	Unfavorable	23%	47%	47%	442	442
	Promising	21%	77%	92%	442	678
	Favorable	56%	95%	95%	442	442
All results are based on 100,000 simulated trials						

50

Adaptation to Gain Efficiency?

- Consider adaptation merely to repower study
 - “We observed a result that was not as good as we had anticipated”

- All GST are within family of adaptive designs
 - Don't we have to be at least as efficient?

- Issues
 - Unspecified adaptations
 - Comparing apples to apples

51

51

Apples with Apples

- Can adapting beat a GST with the same number of analyses?
 - Fixed sample design: $N=1$
 - Most efficient symmetric GST with two analyses
 - $N = 0.5, 1.18$
 - $ASN = 0.6854$
 - Most efficient adaptive design with two possible N
 - $N = 0.5$ and either 1.06 or 1.24
 - $ASN = 0.6831$ (0.34% more efficient)

Table 1: Average and Maximal Sample Sizes of Adaptive Designs in Setting 1

	Number of Continuation Regions							
	1	2	3	4	5	6	7	8
<i>ASN</i>	0.6854	0.6831	0.6828	0.6825	0.6824	0.6824	0.6824	0.6824
<i>% Reduction</i>	Ref	0.34%	0.38%	0.42%	0.43%	0.43%	0.44%	0.44%
<i>Maximal N</i>	1.18	1.24	1.24	1.26	1.26	1.26	1.26	1.28

52

52

Apples with Apples (continued)



- Can adapting beat a GST with the same number of analyses?
 - Fixed sample design: $N=1$
 - Most efficient symmetric GST with two analyses
 - $N = 0.5, 1.18$
 - $ASN = 0.6854$
 - GST with same three analyses
 - $N = 0.5, 1.06$ and 1.24
 - $ASN = 0.6666$ (2.80% more efficient)
 - GST with same five analyses
 - $N = 0.5, 1.01, 1.10, 1.17,$ or 1.31
 - $ASN = 0.6576$ (4.20% more efficient)

53

53

Comments re Conditional Power



- Many propose adaptations based on conditional /predictive power

- Neither have good foundational motivation
 - Frequentists should use Neyman-Pearson paradigm and consider optimal unconditional power across alternatives
 - And conditional/predictive power is not a good indicator in loss of unconditional power
 - Bayesians should use posterior distributions for decisions

- Difficulty understanding conditional / predictive power scales can lead to bad choices for designs

54

54

Comparisons of Designs

- The example used here was a longitudinal study, rather than time to event, though the same issues obtain
- Statistical power
- Sample size accrued
 - With time to event, often all subjects have been accrued when half the statistical information is not yet available
- Calendar time
 - Number of events is more a surrogate for savings in time monitoring subjects and marketing time lost

55

55

Alternative Approaches

Table 1: Comparison of RCT Designs for Example 1

Design	Hypothesized Treatment Effect						
	$\delta = 0$	$\delta = 1.5$	$\delta = 1.6$	$\delta = 1.7$	$\delta = 1.8$	$\delta = 1.9$	$\delta = 2.0$
	Power						
<i>Fxd442</i>	2.5%	55.6%	61.1%	66.3%	71.3%	75.9%	80.0%
<i>Fxd690</i>	2.5%	74.8%	80.0%	84.5%	88.3%	91.4%	93.9%
<i>GST694</i>	2.5%	74.8%	80.0%	84.6%	88.4%	91.4%	93.9%
<i>Adapt</i>	2.5%	60.4%	65.8%	70.8%	75.4%	79.6%	83.4%
<i>Fxd492</i>	2.5%	60.2%	65.8%	71.0%	75.9%	80.2%	84.1%
<i>Fut492</i>	2.5%	59.8%	65.4%	70.6%	75.4%	79.8%	83.7%
<i>OBF492</i>	2.5%	59.6%	65.2%	70.4%	75.3%	79.6%	83.5%
	Expected Number Accrued						
<i>Fxd442</i>	442	442	442	442	442	442	442
<i>Fxd690</i>	690	690	690	690	690	690	690
<i>GST694</i>	694	681	678	675	671	667	662
<i>Adapt</i>	464	496	495	494	492	490	488
<i>Fxd492</i>	492	492	492	492	492	492	492
<i>Fut492</i>	468	488	489	490	490	490	491
<i>OBF492</i>	467	485	485	485	485	484	484

56

56

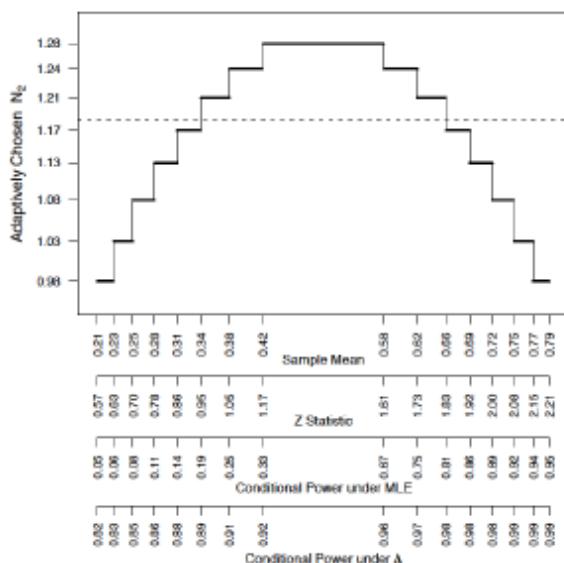
Apparent Problem

- The authors chose extremely inefficient thresholds for conditional power
 - Adaptation region $0.365 < CP_{est} < 0.8$
 - From optimal test, $0.049 < CP_{est} < 0.8$ is optimal
- Of course, we do not always choose the most efficient designs
 - O'Brien-Fleming designs are markedly inefficient for primary endpoint, but do allow adequate sample size for safety and secondary endpoints
- But more careful evaluation can allow us to choose adaptations that satisfy desired operating characteristics

59

59

“Optimal” Adaptive Design



60

60

The Cost of Planning Not to Plan

- Hypothesis testing of a null with fully adaptive trials
 - Statistics: type I error is controlled
 - Game theory: chance of “winning” with completely ineffective therapy is controlled
 - Science:
 - Discrimination of clinically relevant hypothesis may be impaired
 - May be uncertain as to what the treatment has effect on

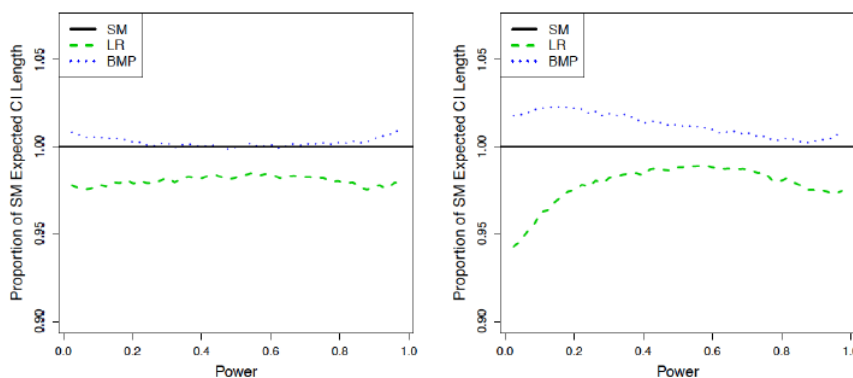
- Frequentist estimation: (Levin, Emerson, Emerson, 2012)
 - Ideally pre-specify the adaptive rule
 - GST methods can be extended to adaptive sampling density
 - When fully adaptive, Brannath, Mehta, Posch (2009) have proposed a very clever method that works reasonably well.

61

61

Comparison of CI Length

- Reference OF design, symmetric (left) or CP-based (right) N_J function, up to 50% increase, $J = 2$



62

62

Proportional Hazards

.....

SSRE with Extreme Treatment Effects

Where am I going?

Design of a RCT is based on a variety of assumptions that may not obtain in practice

Investigators then may have an interest in adjusting the RCT design to better address the actual conditions

65

65

Motivation



- Consider the design of an RCT that investigates prevention strategies in HIV / AIDS
- Our primary clinical endpoint is sero-conversion to HIV positive
- We will randomize individuals 1:1 experimental treatment to control

66

66

Recall



- In the presence of time to event endpoint that is subject to censoring, the most commonly used analyses are the logrank test and the proportional hazards regression model (Cox regression)
- When using PH regression with alternatives that satisfy the PH assumption, statistical information is proportional to the number of events
 - We can separately consider number accrued and calendar time of ending study
- Sample size calculations thus return the number of events that are necessary to obtain desired power
 - There are multiple ways that we can obtain that number of events as a function of
 - Number and timing of accrued subjects
 - Length of follow-up after start of study

67

67

Motivation: HPTN052



- Highly effective treatment and possibly low event rate
- HPTN052: 2011 scientific breakthrough of the year
 - Early vs Delayed ART is effective treatment in the prevention of HIV-1 transmission
 - Design: 188 events anticipated
 - based on (Placebo: 13.2% vs Treatment: 8.3%)
 - Blinded analysis: Total of 28 events
 - Unblinded analysis: 27 from the delayed ART arm
 - HR: 0.04 95% CI 0.01 - 0.27

68

68

Motivation: Partners PrEP

- Highly effective treatment and possibly low event rate
- Partners PrEP: 2012
 - Three arm double-blind trial of daily oral tenofovir (TDF) and emtricitabine/tenofovir (FTC/TDF)
 - 1:1:1 randomization of 4578 serodiscordant couples
 - Study halted 18 months earlier than planned due to demonstrated effectiveness in reduction of HIV-1 transmission
 - Of 78 infections, 18 in tenofovir, 13 in Truvada, 47 in control
 - Reduction in risk of infection 62% (95% CI 34-78%) in tenofovir, 73% (95% CI 49-85%); $p < 0.0001$ vs control
 - Special note: Placebo event rate was 1.99 per 100 PY rather than planned 2.75 per 100 PY

69

69

Motivation: HPTN 083

Article

CLINICAL TRIALS

Evaluating group-sequential non-inferiority clinical trials following interim stopping: The HIV Prevention Trials Network 083 trial

Brett S Hanscom¹, Deborah J Donnell¹, Thomas R Fleming², James P Hughes^{1,2}, Marybeth McCauley³, Beatriz Grinsztejn⁴, Raphael J Landovitz⁵ and Scott S Emerson²

Clinical Trials
1-8
© The Author(s) 2022
Article reuse guidelines:
sagepub.com/journals-permissions
DOI: 10.1177/17407745221118371
journals.sagepub.com/home/ctj
SAGE

Abstract
Background/Aims: The HIV Prevention Trials Network 083 trial was a group-sequential non-inferiority trial designed to compare HIV incidence under a novel experimental regimen for HIV prevention, long-acting injectable cabotegravir, with an active-control regimen of daily oral tenofovir disoproxil fumarate/emtricitabine (brand name Truvada). In March of 2020, just as the trial had completed enrollment, the COVID-19 pandemic threatened to prevent trial participants from attending study visits and obtaining study medication, motivating the study team to update the interim monitoring plan. The Data and Safety Monitoring Board subsequently stopped the trial at the first interim review due to strong early evidence of efficacy.
Methods: Here we describe some unique aspects of the trial's design, monitoring, analysis, and interpretation. We illustrate the importance of computing point estimates, confidence intervals, and p values based on the sampling distribution induced by sequential monitoring.
Results: Accurate analysis, decision-making and interpretation of trial results rely on pre-specification of a stopping boundary, including the scale on which the stopping rule will be implemented, the specific test statistics to be calculated, and how the boundary will be adjusted if the available information fraction at interim review is different from planned. After appropriate adjustment for the sampling distribution and overrun, the HIV Prevention Trials Network 083 trial provided strong evidence that the experimental regimen was superior to the active control.

70

70

Motivation: HPTN 083 Interim Analysis



- DSMB recommends termination at 44 observed events
 - Estimated HR MLE: 0.29 (nominal 95% CI: 0.14–0.58)

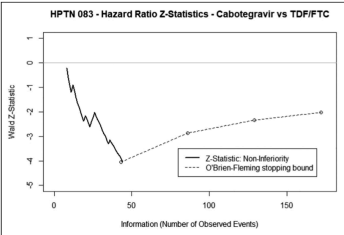


Figure 1. Longitudinal representation of the hazard-ratio z-statistic comparing the experimental product cabotegravir with the active-control product TDF/FTC, just crossing the O'Brien-Fleming monitoring boundary at the first interim analysis.

- Overrun to 52 events at final study analysis
 - Estimated HR MLE: 0.33 (nominal 95% CI: 0.17–0.61)
 - HR Adjusted for GSD: 0.34 (95% CI: 0.18–0.62)

71

71

Motivation: HPTN 083 Interpretation



- “The protocol (or statistical analysis plan) must pre-specify which analytic methods are to be used for computation of the final point estimate, CI, and p value once the trial has stopped, including the choice of adjusted estimator (median unbiased estimator, bias-adjusted mean, etc.), and the outcome-space ordering. In 2019, the Food and Drug Administration convened a panel to address these issues, with resulting guidance including:

“... there are known methods for adjusting estimates to reduce or remove bias associated with adaptations and to improve performance on measures such as the mean squared error.^{1,28} Such methods should be prospectively planned and used for reporting results when they are available.”

72

72

Motivation: HPTN 083 Interpretation



- Changes to GS boundary pre specified by blinded study team
 - Conservatism of OBF boundary meant first analysis stopping threshold establishes superiority as well as noninferiority
- All methods for inference pre specified
 - Ordering based on MLE used for CI and p values
 - (Extreme early conservatism has minimal impact at first analysis)
- Overrun is truly independent of stopping rule
 - Convolution with GS results could be used
 - Most times study teams just update GS boundary to full information
- Trial results establish superiority of cabotegravir

73

73

Issues



- In the first two of these trials the number of events observed was much lower than had been anticipated
- A priori, there are two reasons observed event rates could be lower than anticipated
 - Lower event rate in the control arm that had been guessed
 - Highly effective treatment leads to very few events in the experimental treatment
- In retrospect, both of these trials had both of these problems

74

74

Possible Solutions



- Well-understood methods
 - Wrong baseline event rate
 - Extend planned follow-up time
 - Live with lower power at planned calendar time EOS
 - Adaptive sample size re-estimation based on blinded results
 - Tradeoffs between accrual size and follow-up
 - Highly effective therapy
 - Group sequential design
- Less understood methods
 - Adaptive sample size re-estimation based on blinded results
 - Differentially revise maximum number of events and/or accrual/follow-up based on interim estimates of treatment effect

75

75

Extending Time of Follow-Up



- Under “information time” monitoring, this presents no statistical issues when proportional hazards holds
 - And “information time” monitoring is the usual standard in prespecifying RCT design in the time to event setting, and we would be supposed to do this
- Sometimes, however, we are only willing to believe PH assumption over some shorter time of follow-up
 - National Lung Screening Trial
 - Vaccine trials where need for boosters is not known
- Always, calendar time is ultimately more costly than number of patients
 - Emerson SC, et al. considers tradeoffs between time and number of patients

76

76

Accepting Lower Power



- If the prespecified RCT design defined the maximal statistical information according to calendar time, there is no statistical issue
- Under “information time” monitoring, this represents an unplanned change in the maximal statistical information
 - When this decision is made without knowledge of the unblinded treatment effect, regulatory agencies will usually allow the reporting of a “conditional analysis”
 - But the sponsor will need to be able to convincingly establish that it was still blinded to treatment effect
- Ethics of performing a grossly underpowered study must be considered
 - The predictive value of a “positive” study is greatly reduced

77

77

Blinded Adaptation of Sample Size



- If the prespecified RCT design defined the maximal statistical information according to number of events, then we must be talking about blinded adaptation of accrual size
 - Under PH distribution with PH analysis, no statistical issue
- Under “calendar time” monitoring, this represents an unplanned change in the maximal statistical information
 - When this decision is made without knowledge of the unblinded treatment effect, regulatory agencies will usually allow the reporting of a “conditional analysis”
 - But the sponsor will need to be able to convincingly establish that it was still blinded to treatment effect
- This is likely only credible if you were delaying end of study

78

78

Group Sequential Design

- Instead of a fixed sample design, pre-specify a group sequential design with, say, 10 possible analyses
 - Example: level 0.025, 90% power to detect HR=0.6

```
seqDesign(prob.model = "hazard", alt.hyp = 0.6, nbr.an = 10, power = 0.9)
PROBABILITY MODEL and HYPOTHESES:
Theta is hazard ratio (Treatment : Comparison)
One-sided hypothesis test of a lesser alternative:
  Null hypothesis : Theta >= 1.0      (size = 0.025)
  Alternative hypothesis : Theta <= 0.6 (power = 0.900)
(Emerson & Fleming (1989) symmetric test)
STOPPING BOUNDARIES: Sample Mean scale
```

		Efficacy	Futility
Time 1	(NEv= 17.47)	0.0454	11.8598
Time 2	(NEv= 34.95)	0.2132	2.5280
Time 3	(NEv= 52.42)	0.3568	1.5101
Time 4	(NEv= 69.90)	0.4617	1.1672
Time 5	(NEv= 87.37)	0.5389	1.0000
Time 6	(NEv= 104.85)	0.5974	0.9021
Time 7	(NEv= 122.32)	0.6430	0.8381
Time 8	(NEv= 139.79)	0.6795	0.7931
Time 9	(NEv= 157.27)	0.7093	0.7597
Time 10	(NEv= 174.74)	0.7341	0.7341

79

79

Group Sequential Design

- Stopping boundaries, stopping probabilities

Legend: Fixed (black line), OBFsymm.10 (red line)

Legend: Lower (orange), Upper (green)

80

80

Group Sequential Design

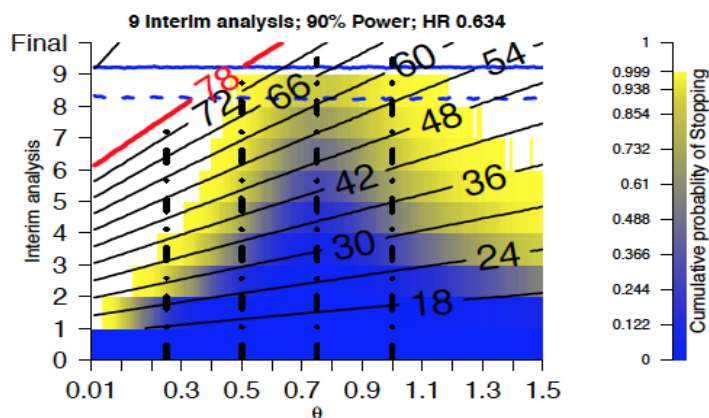
- Using this example, we see that if the true HR was 0.4 or less, we are virtually assured of stopping at the 4th analysis or earlier
- While the maximal number of events was 175, the 4th analysis occurs with 70 events.
- Suppose, a slow accrual of events is due solely to a highly effective treatment
 - Placebo has the planned event rate, Experimental treatment has extremely low event rate
- Relatively frequent monitoring will cause early termination long before the maximal event size needs to be observed
- We examine how calendar time might be affected

81

81

Calendar Time: Half Event Rate

- Stopping probabilities under planned event rate



82

82

Incorporating Lower Event Rates

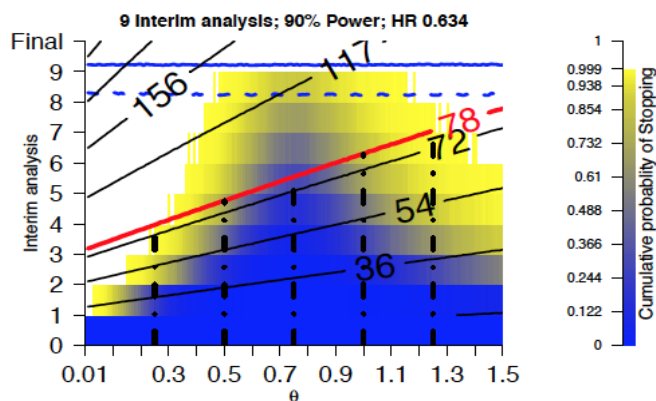
- We have not totally addressed problems that might arise with lower baseline event rates in the control group
 - If the treatment effect is not extreme, then the GSD might dictate that we proceed to the maximal sample size
- One approach is to build in an “escape clause” in the pre-specification of the RCT design
 - “The study will definitely terminate when we have 175 events or at 78 months after start of RCT, whichever comes first.”

83

83

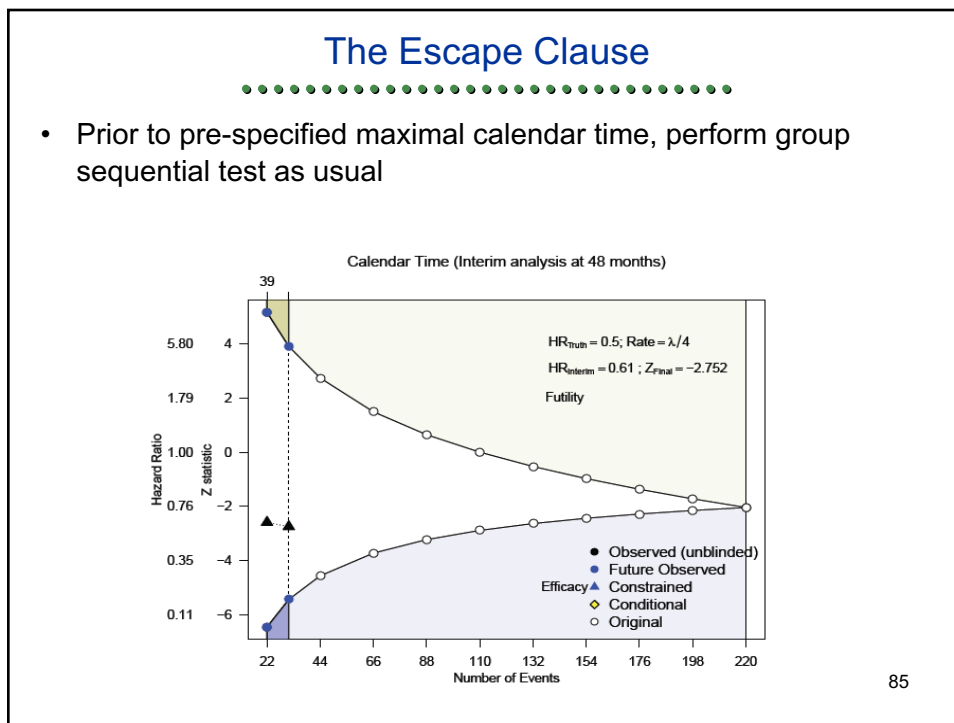
Calendar Time: Half Event Rate

- If control group event rate is halved
 - Power is affected relatively little

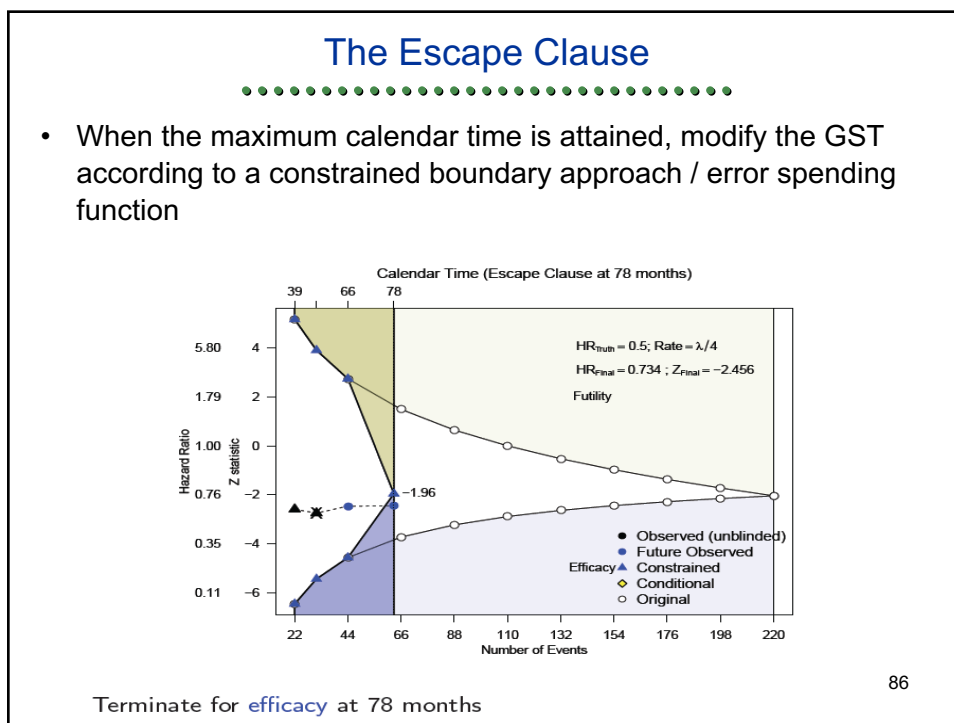


84

84



85



86

Unblinded Adaptation

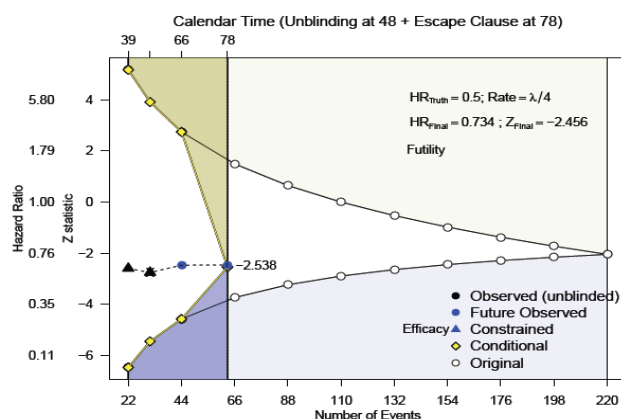
- With unblinded adaptation, we can try to discriminate between
 - Strong treatment effect → choose lower maximal event size
 - Low control event rate → accrue more information
- We will have to decide whether to do adaptation prior to stopping accrual or whether to restart accrual
 - Early adaptation → Less precise estimates of treatment effect
 - Late adaptation → Have to restart accrual

87

87

What if Unblinded?

- When the maximum calendar time is attained, have to adjust the critical value according to the conditional error (CHW) or similar



88

Terminate for **futility** at 78 months (More conservative critical value)

88

Simulations



	HR=0.5 ; $\lambda/4$				HR=0.6343; $\lambda/2$			
	Continue		Restart		Continue		Restart	
	Pres	Cond	Pres	Cond	Pres	Cond	Pres	Cond
1750	68.69	-	68.69	-	67.55	-	67.55	-
3500	90.08	-	80.27	-	88.40	-	79.47	-
Fully Blinded [‡]	90.08	89.72	80.27	76.88	87.61	87.60	79.47	79.51
Avg Rate (80%)	86.33	85.74	78.27	73.91	84.63	84.59	77.55	77.36
Rate Diff (80%)	88.09	86.52	80.27	75.25	86.21	85.69	79.31	78.84
HR (80%)	87.55	86.31	80.10	75.07	86.10	85.58	79.35	78.77

- ▶ GSD (fully blinded procedures) almost efficient to the best *prespecified adaptive design* in context of $\lambda_{\text{Truth}} < \lambda_{\text{Planned}}$
- ▶ However, when integrity of the trial may be compromised and adjustments have to be used (CHW), we lose power
- ▶ The inefficient weighting scheme of CHW results in substantial loss of power particularly with late adaptations.

89

Final Comments



- The group sequential design definitely protects us from the extreme treatment effect
- In general, the group sequential design protected us from problems so long as the event rate was at least 25% of the planned rate
- There was definitely a price to pay when using the adaptive design
 - If the sponsor has access to unblinded results, adjustment for the adaptive analysis must be made
 - There is no allowance for the “escape clause” approach
 - Even more difficulty if non PH is possible

90

90

Nonproportional Hazards

.....

Weighted Logrank Statistics

Where am I going?

Early phase clinical trials sometimes show treatment effects that are more pronounced early or more pronounced late

Weighted versions of the logrank statistic have been proposed to accentuate those portions of the survival curve that are most plausibly different

91

91

Weighted Logrank Statistics

.....

- Choose additional weights to detect anticipated effects

$$W(\beta) = \sum_t w(t) \frac{n_{0t}n_{1t}}{n_{0t} + n_{1t}} \left[\hat{\lambda}_{1t} - e^{\beta} \hat{\lambda}_{0t} \right]$$

$$n_{kt} = N_k \times \Pr(T \geq t, Cens \geq t) \stackrel{ind}{=} N_k S_k(t) \times \Pr(Cens \geq t)$$

$G^{\rho\gamma}$ Family of weighted logrank statistics :

$$w(t) = \left[\hat{S}_\cdot(t) \right]^\rho \left[1 - \hat{S}_\cdot(t) \right]^\gamma$$

92

92

What if No Adjustment?



- Many methods for adaptive designs seem to suggest that there is no need to adjust for the adaptive analysis if there were no changes to the study design
- However, changes to the censoring distribution definitely affect
 - Distribution-free interpretation of the treatment effect parameter
 - Statistical precision of the estimated treatment effect
 - Type 1 error when testing a weak null (e.g., noninferiority)
- Furthermore, “less understood” analysis models prone to inflation of type 1 error when testing a strong null
 - Information growth with weighted log rank tests is not always proportional to the number of events

93

93

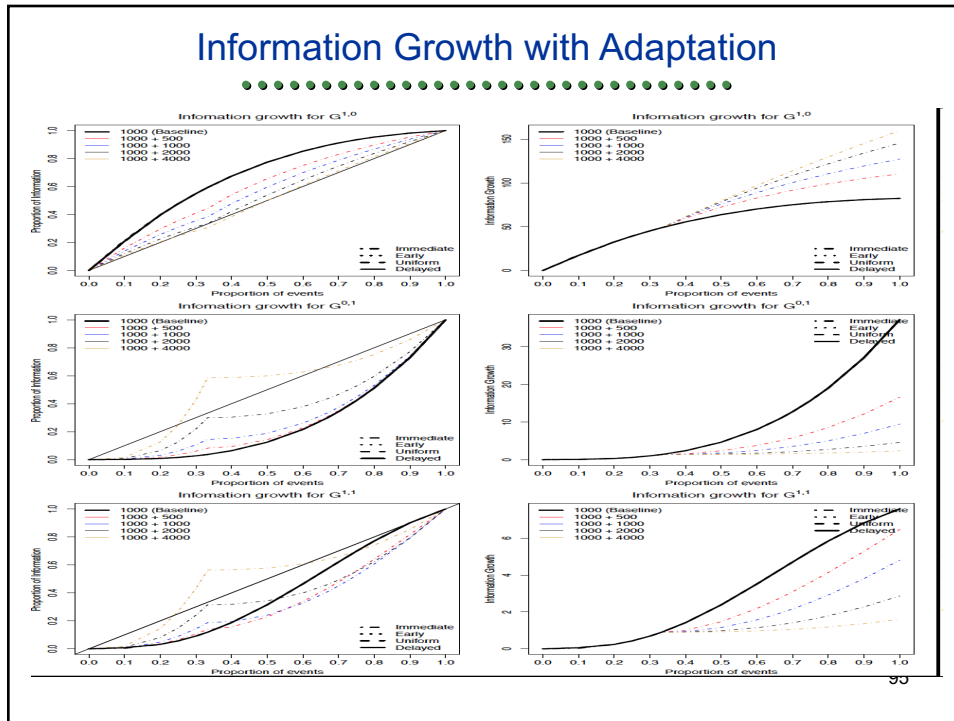
“Intent to Cheat” Zone



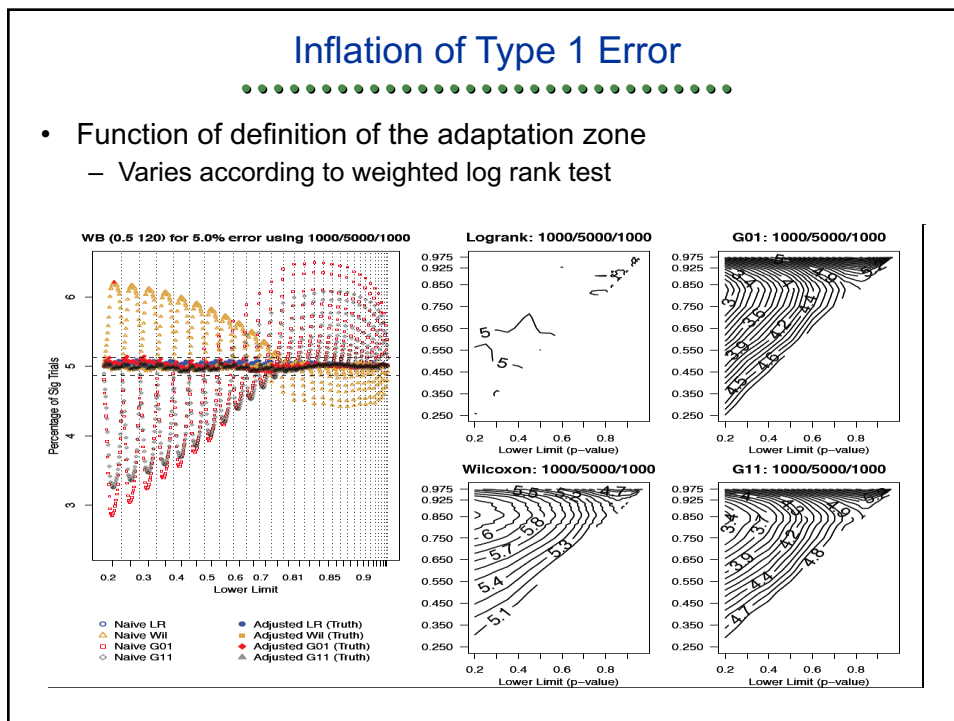
- At interim analysis, choose range of interim estimates that lead to increased accrual of patients
- How bad can we inflate type 1 error when holding number of events constant?
- Logrank test under strong null: Not at all
- Weighted logrank tests: Up to relative increase of 20%
 - Sequela of true information growth
 - **Information growth not linear in number of events**
 - Power largely unaffected, so PPV decreases

94

94



95



96

Comments re WLR



- Hence, unblinded access to trial results can allow an investigator to inflate the type 1 error
- This might not be noticeable to a naïve audience if the number of events stays constant
- Proper handling of information growth can fix this
 - However, description of the information growth is often difficult with weighted log rank statistics

97

97

Nonproportional Hazards



Crossing Survival Curves

Where am I going?

Recently some authors have proposed sequential tests to be used in the presence of crossing survival curves

This example illustrates many of the difficulties inherent in applying time to event analyses

98

98

A Further Example

BIOMETRICS 64, 733–740
 September 2008

DOI: 10.1111/j.1541-0420.2007.00975.x

Comparing Treatments in the Presence of Crossing Survival Curves: An Application to Bone Marrow Transplantation

Brent R. Logan,* John P. Klein, and Mei-Jie Zhang

Division of Biostatistics, Medical College of Wisconsin, 8701 Watertown Plank Road,
 Milwaukee, Wisconsin 53226, U.S.A.

*email: blogan@mcw.edu

SUMMARY. In some clinical studies comparing treatments in terms of their survival curves, researchers may anticipate that the survival curves will cross at some point, leading to interest in a long-term survival comparison. However, simple comparison of the survival curves at a fixed point may be inefficient, and use of a weighted log-rank test may be overly sensitive to early differences in survival. We formulate the problem as one of testing for differences in survival curves after a prespecified time point, and propose a variety of techniques for testing this hypothesis. We study these methods using simulation and illustrate them on a study comparing survival for autologous and allogeneic bone marrow transplants.

KEY WORDS: Censored data; Crossing hazard functions; Generalized linear models; Log-rank test; Pseudo-value approach; Weibull distribution; Weighted Kaplan–Meier statistic.

99

99

Logan, et al.: Motivation

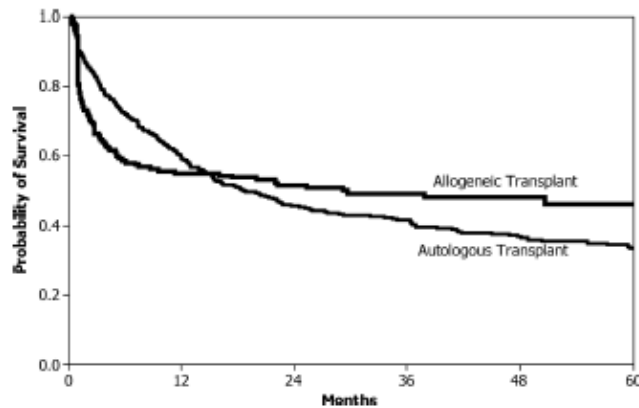


Figure 1. Kaplan–Meier estimate of DFS for follicular lymphoma example, by stem cell source.

100

100

Logan, et al.: Comparisons

- Logrank starting from time 0
- Weighted logrank test ($\rho=0$, $\gamma=1$) from time 0
- Survival at a single time point after time t_0
- Logrank starting from time t_0
- Weighted area between survival curves (restricted mean)
 - Most weight after time t_0
- Pseudovalues after time t_0
- Combination tests (linear and quadratic)
 - Compare survival at time t_0
 - Compare hazard ratio after time t_0

101

101

Logan, et al.: Simulations

Comparing Treatments with Crossing Survival Curves

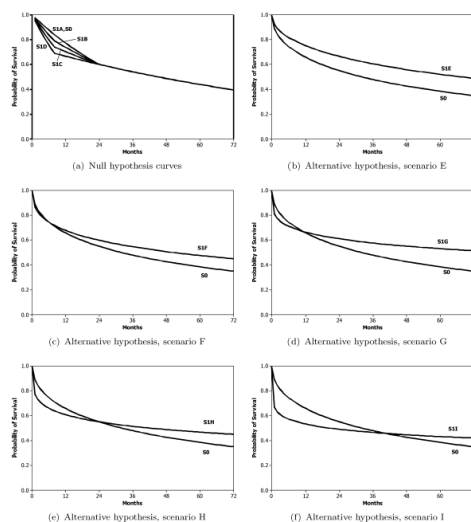


Figure 2. Survival curves for treatment (S1) and control (S0) groups used in simulations. Curves for the null hypothesis simulations are shown in (a) for each of the four scenarios, and curves for the alternative hypothesis simulations are shown in (b)–(f) for the four scenarios.

102

102

Logan, et al.: Results



Table 2

Average rejection rates for 11 tests adjusted using ANOVA for censoring pattern. Rejection rates given by scenario using model (12). The last two rows refer to the log-rank (LR) test and weighted log-rank (WLR) tests starting at time 0. $t_0 = 24$.

Method	Equation	Scenario				
		E	F	G	H	I
$Z_{CLL}(24)$	(1)	62.4	15.3	21.1	4.7	21.8
$Z_{CLL}(48)$	(1)	70.1	32.9	65.1	21.5	6.8
$Z_{CLL}(72)$	(1)	71.2	44.5	85.1	46.1	25.9
$Z_{WKM}(t_0)$	(2)	75.8	35.0	66.3	20.3	6.0
$\chi^2_{PSV}(t_0)$	(3)	74.8	32.0	61.2	16.4	4.8
$Z_{LR}(t_0)$	(4)	30.7	36.5	85.4	71.7	82.6
$Z_{OLS}(t_0)$	(5)	74.7	43.9	84.1	43.4	23.6
$Z_{SP,P}(t_0)$	(6)	76.9	40.2	74.8	29.6	10.7
$\chi^2(t_0)$	(7)	67.2	36.7	83.1	61.1	81.0
Log rank		78.0	28.9	47.0	8.6	22.2
Weighted log rank $\rho = 0, \gamma = 1$		64.7	49.7	93.8	70.0	64.6

103

103

Logan, et al.: Critique



- In considering the combination tests, crossing survival curves might have
 - No difference at time t_0 (perhaps we are looking for equivalence)
 - Higher hazard after time t_0
- Presumably, the authors are interested in the curve that is higher at longer times post treatment
 - The authors did not describe how to use their test in a one-sided setting
- **PROBLEM:** The authors do not seem to be considering the difference between crossing survival curves and crossing hazard functions
 - Higher hazard over some period of time does not imply lower survival curves

104

104

Relevance to Today

- Even experts in survival analysis sometimes lose track of the way that time to event analyses behave, relative to our true goals

Lifetime Data Anal (2015) 21:218–240
 DOI 10.1007/s10985-014-9298-4

Group sequential tests for long-term survival comparisons

Brent R. Logan · Shuyuan Mo

Received: 6 February 2014 / Accepted: 2 July 2014 / Published online: 23 July 2014
 © Springer Science+Business Media New York 2014

Abstract Sometimes in clinical trials, the hazard rates are anticipated to be nonproportional, resulting in potentially crossing survival curves. In these cases, researchers are usually interested in which treatment has better long-term survival. The log-rank test and the weighted log-rank test may not be appropriate or efficient to use here, because they are sensitive to differences in survival at any time and don't just focus on long-term outcomes. Also in a prospective clinical trial, patients are entered sequentially over calendar time, so that group sequential designs may be considered for ethical, administrative and economic concerns. Here we develop group sequential methods for testing the null hypothesis that the survival curves are identical after a prespecified time point. Several classes of tests are considered, including an integrated difference in survival probabilities after this time point, and linear or quadratic combinations of two component test statistics (pointwise comparisons of survival at the time point and comparisons of hazard rates after the time point). We examine the type I errors, stopping probabilities, and powers of these tests through simulation studies under the null and different alternatives, and we apply them to a real bone marrow transplant clinical trial.

Keywords Crossing hazards · Crossing survival curves · Late survival difference · Group sequential test · Error-spending methods

107

107

Proportional Hazards

- All statistics behave in sensible fashion with increasing time

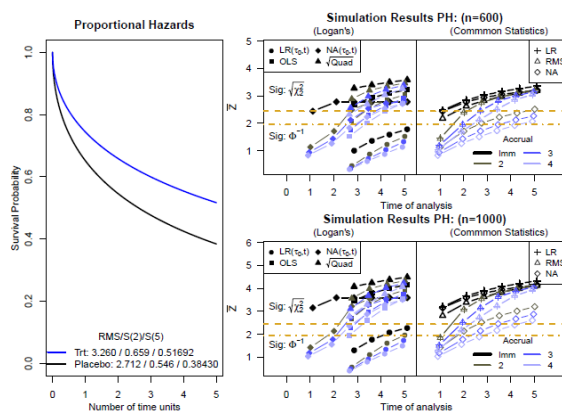


Figure 7.2: Survival curves under proportional hazards alternatives for $n=600$ and 1000 . A positive Z corresponds to the treatment being superior relative to the placebo. Even with censoring, the alternatives are consistently positive (respectively in quadrant I) for the commonly used or composite statistics.

108

108

Crossing Survival Curves

- As expected, early analyses show evidence opposite to later analyses

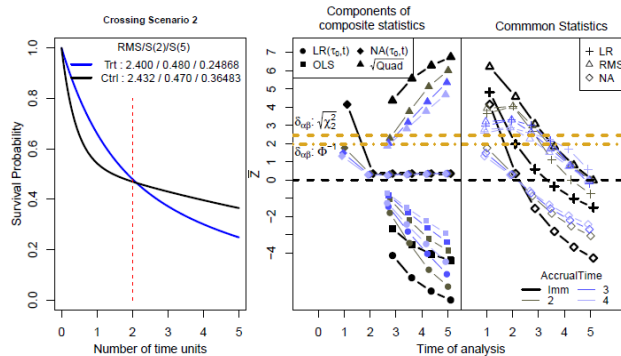


Figure 7.4: Simulated crossing hazards, crossing survival curves where we added variations to whether the curves crosses just before 2, and after 2. The combination of composite alternatives changes from Quadrant III (-,-), Quadrant III/IV(0, -), Quadrant IV (+,-) with the net result providing conclusion of preferring the placebo over treatment. The commonly used statistics is seen to have a changing alternative that switches from positive to negative.

109

109

Stochastically Ordered w/ Crossing Hazards

- Proposed composite statistics disagree with clinical judgment

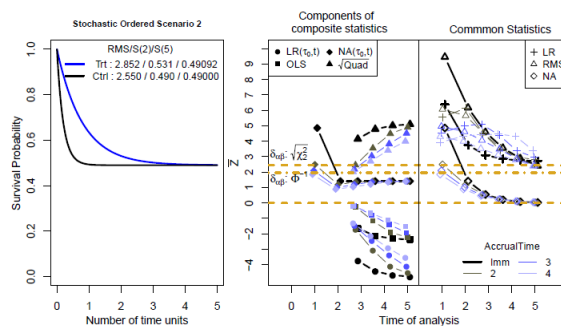


Figure 7.3: Standardized alternative at different interim analyses under stochastic ordered, crossing hazards survival curves with various accrual patterns for scenario 2. Survival curves for scenario 2 under stochastic ordering with roughly 50% probability of survival. The combination of alternatives for NA(τ_0, t) and LR(τ_0, t) resides in quadrant IV and describes conflicting conclusions prior to crossing and after crossing. This observation is as speculated and described in Table 7.1. The net effect for the linear composite statistics conclude placebo is better than the treatment despite treatment being better at all times from 0 to 5 relative to placebo.

110

110

Proportional Hazards



Availability of Surrogate Data

Where am I going?

Methods for preserving type 1 errors presume an accurate representation of the statistical information available at the adaptive analysis

With time to event data (as well as other longitudinal endpoints), however, we may have information on surrogate prognostic endpoints.

To the extent that those surrogate endpoints inform the adaptation of the clinical trial, we may not be adequately preserving the type 1 error

111

111

Special Issues



- A basic premise of adaptive methods is that we can control the type 1 error, even when we have re-designed the trial based on interim estimates of the treatment effect
- Two special scenarios that we need to examine more closely
 - Do the interim statistics used in adjusting critical values truly contain all the information we had at our disposal?
 - Have we quantified the information growth correctly when using those statistics?

112

112

Protecting Type I Error

.....

- Test based on weighted averages of incremental test statistics
 - Allow arbitrary weights W_j specified by stage $j-1$

$$Z = \frac{\sum_{k=1}^J \sqrt{W_k} Z_k^*}{\sqrt{\sum_{k=1}^J W_k}} \quad \bigcap_{k=1}^J H_{0,j} \sim N(0,1)$$

$$Z = \frac{\sum_{k=1}^J \sqrt{W_k} \Phi^{-1}(1 - P_k^*)}{\sqrt{\sum_{k=1}^J W_k}} \quad \bigcap_{k=1}^J H_{0,j} \sim N(0,1)$$

113

113

Complications: Longitudinal Outcomes

.....

- Bauer and Posch (2004) noted that in the presence of incomplete data, partially observed outcome data may be informative of the later contributions to test statistics
 - E.g., tumor progression and overall survival
- This can be a large problem if we allow adaptation to a much smaller sample size
 - Data quite often becomes available between database lock and a DSMB meeting

114

114

Complications: Longitudinal Outcomes

.....

- We need to make distinctions between
 - Independent subjects accrued at different stages
 - Statistical information about the primary outcome available at different analyses

- Owing to delayed observations, contributions to the primary test statistic at the *k*-th stage may come from subjects accrued at prior stages
 - Some information is typically available between “database lock” and DSMB meeting
 - Baseline and secondary outcome data available at prior analyses on censored subjects may inform the value of future data

115

115

Data at *j*-th Analysis: Delayed Outcome

.....

- Subjects accrued at different stages are independent
- Some data is “missing”

At <i>k</i> th interim analysis	Incremental	Cumulative
Sample size (stat info)	N_k^*	$N_k = N_1^* + \dots + N_k^*$
Baseline data	\bar{X}_k^*	$\bar{X}_k = (\bar{X}_1^*, \dots, \bar{X}_k^*)$
1° outcome data (msgng, observed)	$\bar{Y}_k^{*M}, \bar{Y}_k^{*O}$	\bar{Y}_k^M, \bar{Y}_k^O
2° outcome data	\bar{W}_k^*	$\bar{W}_k = (\bar{W}_1^*, \dots, \bar{W}_k^*)$
Estimated treatment effect	$\hat{\theta}_k^* = \hat{\theta}_k^*(N_k^*, \bar{X}_k^*, \bar{Y}_k^{*O}, \bar{Y}_k^{*M})$	$\hat{\theta}_k = \frac{\sum_{j=1}^k N_j^* \hat{\theta}_j^*}{N_k}$
Normalized Z statistic	Z_k^*	$Z_k = \frac{\sum_{j=1}^k \sqrt{N_j^*} Z_j^*}{\sqrt{N_k}}$
Fixed sample P value	P_k^*	

116

116

Major Problem: Delayed Outcome

- When sample size N_j^* and parameter θ_j adaptively chosen based on data from prior stages $1, \dots, j-1$, some aspect of the “future” contributions may already be known

At k th interim analysis	Incremental	Cumulative
Sample size	$N_k^* = N_k^*(N_{k-1}, \bar{X}_{k-1}, \bar{W}_{k-1}, \bar{Y}_{k-1}^{*O}, \bar{Y}_{k-2}^{*M})$	N_k
Estimated treatment effect	$\hat{\theta}_k^* = \hat{\theta}_k^*(N_k^*, \bar{X}_k^*, \bar{Y}_k^{*O}, \bar{Y}_{k-1}^{*M})$	$\hat{\theta}_k^* = \frac{\sum_{j=1}^k N_j^* \hat{\theta}_j^*}{N_k}$

Impact: (One statistician's mean is another statistician's variance)

$$\text{corr}(\bar{Y}_k^{*M}, \bar{W}_k^*) \neq 0 \text{ or } \text{corr}(\bar{Y}_k^{*M}, \bar{X}_k^*) \neq 0 \Rightarrow \hat{\theta}_k^* | N_k^* \text{ not indep of } \hat{\theta}_{k+1}^* | N_{k+1}^*$$

$\hat{\theta}_k^* | N_k^*$ is potentially biased for θ_k and not approximately normal

117

117

Potential Solutions

- Jenkins, Stone & Jennison (2010)
 - Only use data available at the k -th stage analysis
 - (Analogous to phase 2 followed by phase3)
- Irle & Schaefer (2012)
 - Prespecify how the full k -th stage data will eventually contribute to the estimate of θ_k
 - (Pretends that we already had that information when adapting)
- Magirr, Jaki, Koenig & Posch (2014, arXiv.org)
 - Assume worst case of full knowledge of future data and sponsor selection of most favorable P value
 - (Adds multiple comparison issues to the Irle & Schaefer approach)

118

118

Comments: Burden of Proof Dilemma



- There is a contradiction of standard practices when viewing the incomplete data
 - We would never accept the secondary outcomes as validated surrogates
 - But we feel that we must allow for the possibility that the secondary outcomes were perfectly predictive of the eventual data
- We are in some sense preferring mini-max optimality criteria over a Bayes estimator

119

119

Comments: Impact on RCT Design



- The candidate approaches will protect the type 1 error, but the impact on power (and PPV) is as yet unclear
- Weighted statistics are not based on minimal sufficient statistics
 - But greatest loss in efficiency comes from late occurring adaptive analyses with large increases in maximal statistical information
 - Time to event will not generally have this, though might with the Irlle & Schafer approach
- The adaptation is based on imprecise estimates of the estimates that will eventually contribute to inference
- We may have to eventually either
 - Ignore some observed data (JS&S, I&S), or
 - Adjust for worst case multiple comparisons

120

120

Remaining Questions 1



- How much inflation of the type 1 error is possible due to final data available between database lock and DSMB meeting?
- Varies according to
 - Timing of adaptive analysis,
 - Amount of existing information not included in the interim analysis, and
 - Adaptive rule
- For Proschan & Hunsberger type modifications at $n_1 = 0.5$
 - 25-40% relative increase in type 1 error plausible

121

121

Remaining Questions 2



- Can we treat surrogate information as predictors in a missing data problem?
 - Missing data due to administrative censoring is MAR
 - Impute a value of Z_1
 - If we presume any “imputation” is based on all informative data, then we remove the issue of surrogacy
 - A very strong presumption
 - However,
 - We have to estimate the information growth for our imputed Z statistic for the interim treatment effect estimate
 - The final statistic will not be a weighted sum of that imputed value, so we will introduce noise into our model
- Preliminary results: The loss of precision from imputing the missing data is too great when used in the CHW adjustment

122

122

Remaining Questions 3



- If we have to adjust for all the information eventually coming from the stage 1 data
 - What is the impact of using imprecise data in the adaptive rule?
 - What is the impact of adjusting the analysis for the correct data?
- Very preliminary results:
 - With a fairly efficient adaptive rule:
 - The efficiency gained from an efficient adaptive rule is fairly small.
 - If we use a relatively efficient adaptive rule, then adjustment for the adaptation has negligible effect on inference
 - Adaptive rules that are less efficient lead to greater differences, but it is hard to identify the part that comes from a bad adaptive rule and adjustment for data not observed at the time of adaptation.

123

123

Final Comments



- There is still much for us to understand about the implementation of adaptive designs
- Most often the “less well understood” part is how they interact with particular data analysis methods
 - In particular, the analysis of censored time to event data has many scientific and statistical issues
- How much detail about accrual patterns, etc. do we want to have to examine for each RCT?
- How much do we truly gain from the adaptive designs?
 - (Wouldn't it be nice if statistical researchers started evaluating their new methods in a manner similar to evaluation of new drugs?) ¹²⁴

124

Bottom Line



- There is no substitute for planning a study in advance
 - At Phase 2, adaptive designs are clearly useful to better control parameters leading to Phase 3
 - Most importantly, learn to take “NO” for an answer
 - At Phase 3, it is less clear whether much is gained from unblinded adaptation
 - And scientific / statistical credibility can suffer

- **“Opportunity is missed by most people because it is dressed in overalls and looks like work.”** -- Thomas Edison

- In clinical science, it is the steady, incremental steps that are likely to have the greatest impact.

125

125