

:

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 23:

Case Study: Gram Negative Sepsis

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

1

Case Study

.....

A RCT of an Antibody to Endotoxin in the Treatment of Gram Negative Sepsis

2

:

Statistical Design

- Steps
 - Defining the probability model
 - Defining the comparison group, primary endpoint, analysis model
 - Defining the statistical hypotheses
 - Null, alternative
 - Defining the statistical criteria for evidence
 - Type I error, power
 - Determining the sample size
 - At each analysis, and maximal sample size
 - Evaluating the operating characteristics
 - Planning for monitoring
 - Updating stopping boundaries according to actual conditions
 - Plans for analysis and reporting results
 - Inference adjusted for sequential sampling plan

3

3

Chosen Endpoint / Study Structure

- Primary endpoint: 28 day all cause mortality
- Comparison group
 - Randomized, double blind, placebo controlled
 - 1:1 randomization
- Hypotheses
 - Treatment effect of 7% absolute increase in mortality
 - On placebo: 30% 28 day mortality
 - On antibody: 23% 28 day mortality
- Burden of proof: Registrational trial
 - Type I error: one-sided 0.025
 - Statistical power: high

4

4

:

Creating a RCT Design : Model



- `seqDesign (`
 `prob.model = "proportions",` (vs several others)
 `arms = 2,` (default)
 `ratio = 1,` (default)
 `variance = "alternative",` (default vs "null","intermediate", number)
 `(many others)`
 `)`

5

5

Creating a RCT Design : Hypotheses



- `seqDesign (`
 `null.hypothesis = 0.30,`
 `alt.hypothesis = 0.23,` (vs unspecified)
 `variance = "alternative",` (vs "null","intermediate" or number)
 `(many others)`
 `)`

6

6

:

Creating a RCT Design : Fixed Sample



- `seqDesign (`
 `nbr.analyses = 1,` (default)
 `test.type = "less",` (vs "greater","two.sided")
 `size = 0.025,` (default)
 `(many others)`
 `)`

7

7

Creating a RCT Design : Statistical Task



- `seqDesign (`
 `alt.hypothesis = 0.23,` (vs unspecified)
 `power = "calculate",` (vs a number, say, 0.90)
 `sample.size = 1700,` (vs unspecified)
 `(many others)`
 `)`

8

8

Creating a RCT Design : Output

- seqDesign (
 - display.scale = "X", (default vs any boundary scale)
 - (many others)

9

9

Creating a RCT Design : GUI via Rcmdr (soon)

The screenshot shows the R Commander interface. On the left, the 'Update Design' dialog box is open, showing various parameters for a 2-Arm Test of Proportions. The 'Test Type' is set to 'less', 'Variance Type' is 'null', 'Significance Level' is 0.025, and 'Power' is 0.9. The 'Script Window' on the right contains the following R code:

```
design.3 <- seqDesign( prob.model='proportions', arms=2,
  null.hypothesis=0.3, alt.hypothesis=0.23, test.type='less',
  variance='alternative', alpha=0.025, power=0.90, nbr.analyses=1, ratio=c(1),
  early.stopping='null', F=c(1), R=c(0), A=c(0))
design.3
```

The 'Output Window' shows the execution of this code, displaying the call to seqDesign and the resulting design object. The 'Messages' window at the bottom shows a note about the R Commander version.

10

:

Evaluation of Fixed Sample Designs



- Clinical trial design is most often iterative
 - Specify an initial design
 - Evaluate operating characteristics
 - Modify the design
 - Iterate

11

11

Fixed Sample Operating Characteristics



- Level of Significance (often pre-specified)
- Sample size requirements
 - Scientific / regulatory credibility; Feasibility
- Power Curve
 - Under null (type I error); design alternative
- Decision Boundary
 - Clinical significance
- Frequentist / Bayesian inference on the Boundary
 - Clinical significance of hypotheses discriminated
 - Sensitivity to Bayesian priors

12

12

:

Evaluation in RCTdesign

- Printing
 - Design: `changeSeqScale (x, scale)`
 - Operating characteristics: `seqOC (x, theta, power)`
 - Inference at boundary: `seqInference (x, theta, power)`
 - Everything: `seqEvaluate (x, theta, power)`
- Plotting
 - `plot (x); seqPlotBoundary (x)`
 - `seqPlotASN (x)`
 - `seqPlotPower (x)`
 - `seqPlotStopProb (x)`
 - `seqPlotInference (x)`

13

13

RCTdesign: Creation of Design

```
> fxd90 <- seqDesign("proportions", null=0.30, alt=0.23, power=0.90)
> fxd90
```

Call:

```
seqDesign(prob.model = "proportions", null.hypothesis = 0.3,
          alt.hypothesis = 0.23, power = 0.9)
```

PROBABILITY MODEL and HYPOTHESES:

Theta is difference in probabilities (Treatment - Comparison)

One-sided hypothesis test of a lesser alternative:

Null hypothesis : $\Theta \geq 0.00$ (size = 0.025)Alternative hypothesis : $\Theta \leq -0.07$ (power = 0.900)

(Fixed sample test)

STOPPING BOUNDARIES: Sample Mean scale

Efficacy Futility

Time 1 (N= 1660.17) -0.0423 -0.0423

14

14

:

```

RCTdesign: evalGST
.....
> seqEvaluate(fxd90)

Stopping Boundaries:
  Anlys SampSize CrudeEst      Z  FxdP Hnoninf
Eff     1 1660.173 -0.0423 -1.96 0.025      NA
Fut     1 1660.173 -0.0423 -1.96 0.025      NA

ASN and Cumulative Stopping Probability at Each Analysis
Power TrueEff AvgSampSiz CumStpPrb 1
0.975 -0.0847 1660.173          1
0.950 -0.0778 1660.173          1
0.900 -0.0700 1660.173          1
0.800 -0.0605 1660.173          1

Inference at the Stopping Boundaries
  Anlys SampSize      BAM  CIlo.m CIhi.m Pval.m
Eff     1 1660.173 -0.0423 -0.0847      0 0.025
Fut     1 1660.173 -0.0423 -0.0847      0 0.025

```

15

Alternative Fixed Sample Designs

.....

- Alternative sample sizes
- Sensitivity to assumptions about variability
 - Comparison of means, geometric means
 - Need to estimate variability of observations
 - Comparison of proportions, odds, rates
 - Need to estimate event rate
 - Comparison of hazards
 - Need to estimate number of subjects and time required to observe required number of events

16

:

Group Sequential Stopping Rules

.....

19

Stopping Rules

.....

- **Basic Strategy**
 - Find stopping boundaries at each analysis such that desired operating characteristics (e.g., type I and type II statistical errors) are attained
- **Issues**
 - Conditions under which the trial might be stopped early
 - When to perform analyses
 - Test statistic to use
 - Relative position of boundaries at successive analyses
 - Desired operating characteristics

20

20

:

Boundary Scales



- Stopping boundaries can be defined on a wide variety of scales
 - Sum of observations
 - Point estimate of treatment effect
 - Normalized (Z) statistic
 - Fixed sample P value
 - Error spending function
 - Conditional probability
 - Predictive probability
 - Bayesian posterior probability

21

21

Creating a RCT Design : Sequential



- `seqDesign (`
 - `nbr.analyses = 4,` (default is 1)
 - `sample.size = (1:4)/4*1700,` (default spacing)
 - `design.family = "X",` (default)
 - `early.stopping = "both",` (default vs "null", "alternative")
 - `P = c(1,1),` (default corresponds to OBF)
 - `A = c(0,0),` (default)
 - `R = c(0,0),` (default)
 - `minimum.constraint = ,` (default)
 - `maximum.constraint = ,` (default)
 - `exact.constraint = ,` (default)
 - (many others)

22

22

:

Evaluation of Designs



- Process of choosing a trial design
 - Define candidate design
 - Evaluate operating characteristics
 - Modify design
 - Iterate

23

23

Evaluation of Designs: Fixed Sample



- Operating characteristics for fixed sample studies
 - Level of Significance (often pre-specified)
 - Sample size requirements
 - Power Curve
 - Decision Boundary
 - Frequentist inference on the Boundary
 - Bayesian posterior probabilities

24

24

Evaluation of Designs: Sequential

- Additional operating characteristics for group sequential studies
 - Probability distribution for sample size
 - Stopping probabilities
 - Boundaries at each analysis
 - Frequentist inference at each analysis
 - Point estimates: Bias adjusted, Median unbiased, MLE, UMVUE
 - CI: LR, MLE, stagewise orderings for GSD; BMP for adaptive
 - P values: LR, MLE, stagewise orderings for GSD; BMP for adaptive
 - Bayesian inference at each analysis
 - Conjugate normal priors
 - Futility measures at each analysis
 - Conditional power for arbitrary hypotheses
 - Predictive power for arbitrary conjugate normal priors

25

25

O'Brien-Fleming Symmetric

```
> obf <- seqDesign("prop",null=0.30,test.type="less",
+ sample.size=1700,nbr.analyses=4,P=1,power=0.9)
> seqEvaluate(obf)
```

Stopping Boundaries:

	Anlys	SampSize	CrudeEst	Z	FxdP	Hnoninf
Eff	1	425	-0.1709	-4.0065	0.0000	0.9997
Eff	2	850	-0.0855	-2.8330	0.0023	0.9774
Eff	3	1275	-0.0570	-2.3131	0.0104	0.8763
Eff	4	1700	-0.0427	-2.0032	0.0226	NA
Fut	1	425	0.0855	2.0032	0.9774	0.0003
Fut	2	850	0.0000	0.0000	0.5000	0.0226
Fut	3	1275	-0.0285	-1.1566	0.1237	0.1237
Fut	4	1700	-0.0427	-2.0032	0.0226	NA

ASN and Cumulative Stopping Probability at Each Analysis under Alternatives

Power	TrueEff	AvgSampSiz	CumStpPrb 1	CumStpPrb 2	CumStpPrb 3	CumStpPrb 4
0.975	-0.0855	1098.676	0.0226	0.5026	0.8897	1
0.950	-0.0786	1162.491	0.0153	0.4144	0.8351	1
0.900	-0.0706	1236.314	0.0095	0.3213	0.7603	1
0.800	-0.0610	1315.958	0.0053	0.2309	0.6675	1

Inference at the Stopping Boundaries

	Anlys	SampSize	BAM	CIlo.m	CIhi.m	Pval.m
Eff	1	425	-0.1624	-0.2242	-0.0866	0.0000
Eff	2	850	-0.0795	-0.1296	-0.0250	0.0024
Eff	3	1275	-0.0543	-0.0957	-0.0068	0.0123
Eff	4	1700	-0.0427	-0.0855	0.0000	0.0250
Fut	1	425	0.0770	0.0011	0.1387	0.9765
Fut	2	850	-0.0060	-0.0605	0.0442	0.4011
Fut	3	1275	-0.0312	-0.0786	0.0102	0.0672

26

26

:

O'Brien-Fleming Efficacy with Futility

```
> fut <- update(obf, P=c(1,0.8))
> seqEvaluate(fut)
```

Stopping Boundaries:

	Anlys	SampSize	CrudeEst	Z	FxdP	Hnoninf
Eff	1	425	-0.1695	-3.9756	0.0000	0.9997
Eff	2	850	-0.0848	-2.8112	0.0025	0.9766
Eff	3	1275	-0.0565	-2.2953	0.0109	0.8744
Eff	4	1700	-0.0424	-1.9878	0.0234	NA
Fut	1	425	0.0473	1.1082	0.8661	0.0076
Fut	2	850	-0.0097	-0.3211	0.3741	0.0625
Fut	3	1275	-0.0310	-1.2577	0.1043	0.1768
Fut	4	1700	-0.0424	-1.9878	0.0234	NA

ASN and Cumulative Stopping Probability at Each Analysis under Alternatives

Power	TrueEff	AvgSampSiz	CumStpPrb 1	CumStpPrb 2	CumStpPrb 3	CumStpPrb 4
0.975	-0.0865	1079.055	0.0266	0.5292	0.9052	1
0.950	-0.0794	1140.971	0.0188	0.4409	0.8556	1
0.900	-0.0713	1211.075	0.0133	0.3495	0.7875	1
0.800	-0.0615	1283.396	0.0110	0.2654	0.7038	1

Inference at the Stopping Boundaries

	Anlys	SampSize	BAM	CIlo.m	CIhi.m	Pval.m
Eff	1	425	-0.1610	-0.2228	-0.0852	0.0000
Eff	2	850	-0.0791	-0.1289	-0.0243	0.0026
Eff	3	1275	-0.0548	-0.0955	-0.0064	0.0129
Eff	4	1700	-0.0437	-0.0865	0.0000	0.0250
Fut	1	425	0.0378	-0.0371	0.1005	0.8458
Fut	2	850	-0.0173	-0.0707	0.0341	0.2628
Fut	3	1275	-0.0348	-0.0821	0.0076	0.0530
Fut	4	1700	-0.0437	-0.0865	0.0000	0.0250

21

27

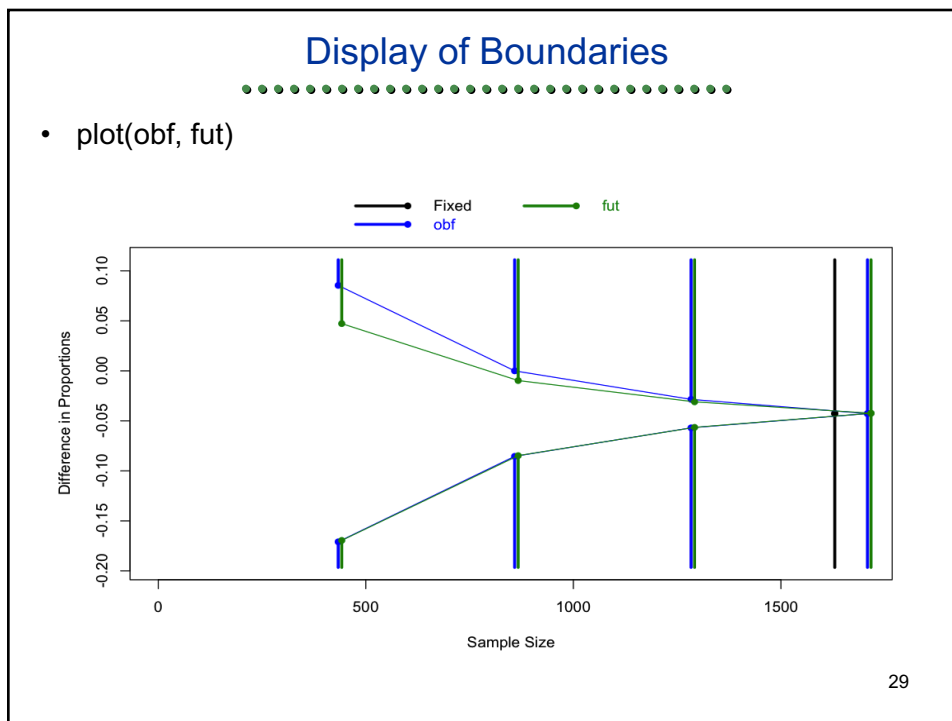
Plotting Stopping Boundaries

- `plot()` or `seqPlotBoundary()`
 - Arbitrary number of designs
 - By default, include fixed design with same power as first specified design
- Axes
 - Y-axis: Critical values for stopping on arbitrary boundary scale
 - MLE scale as default
 - X-axis: Statistical information
 - Sample size (number of events for hazards probability models)
- Display of boundaries
 - Vertical lines are true stopping regions
 - Critical values connected to allow better visualization of boundary shape

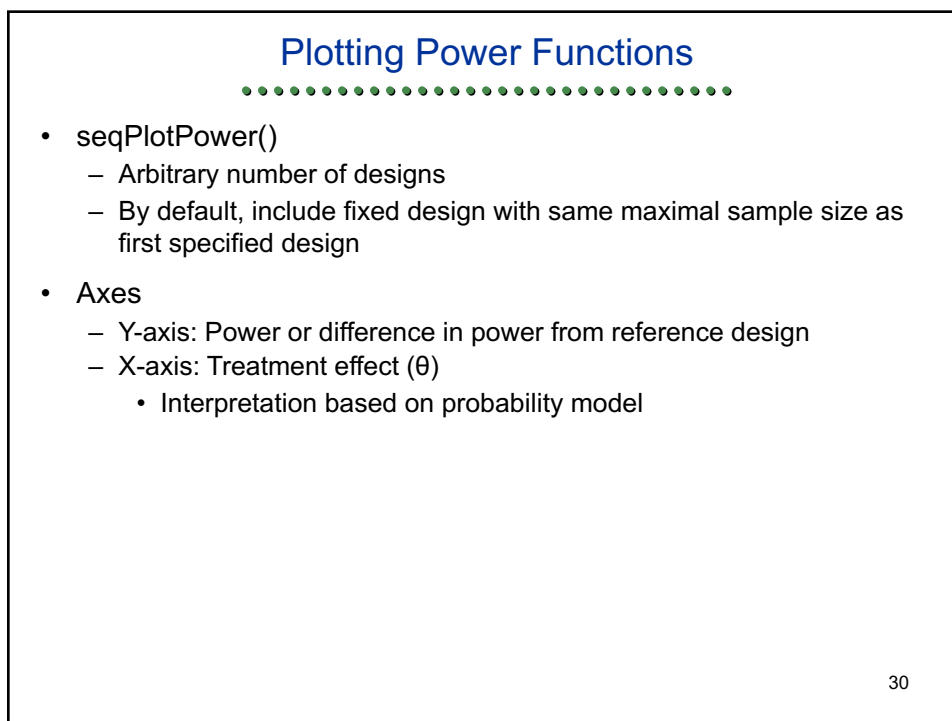
28

28

:

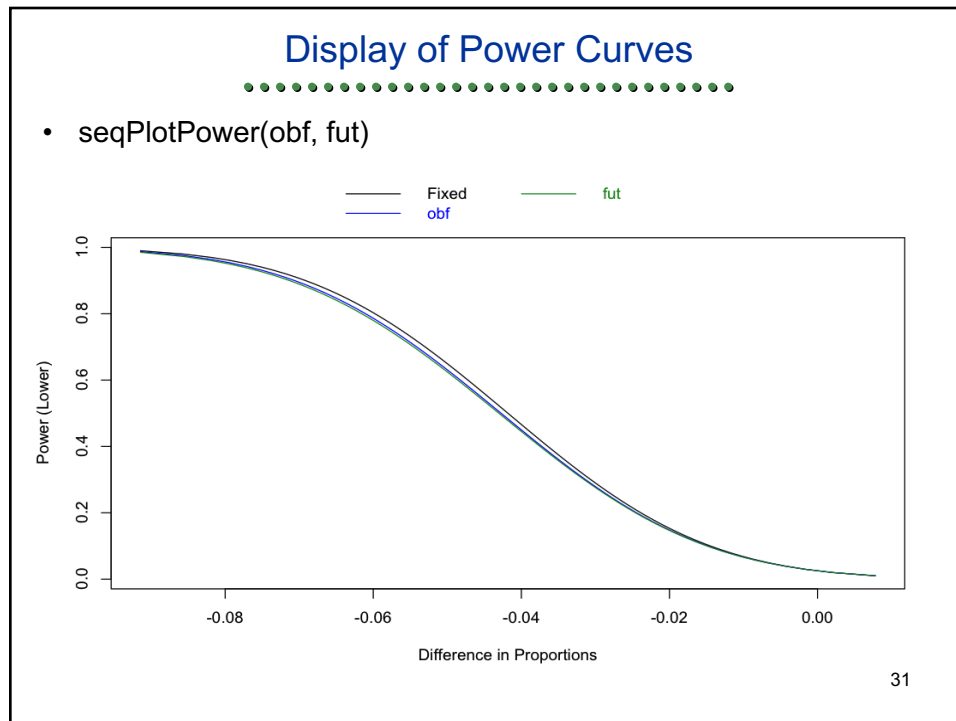


29

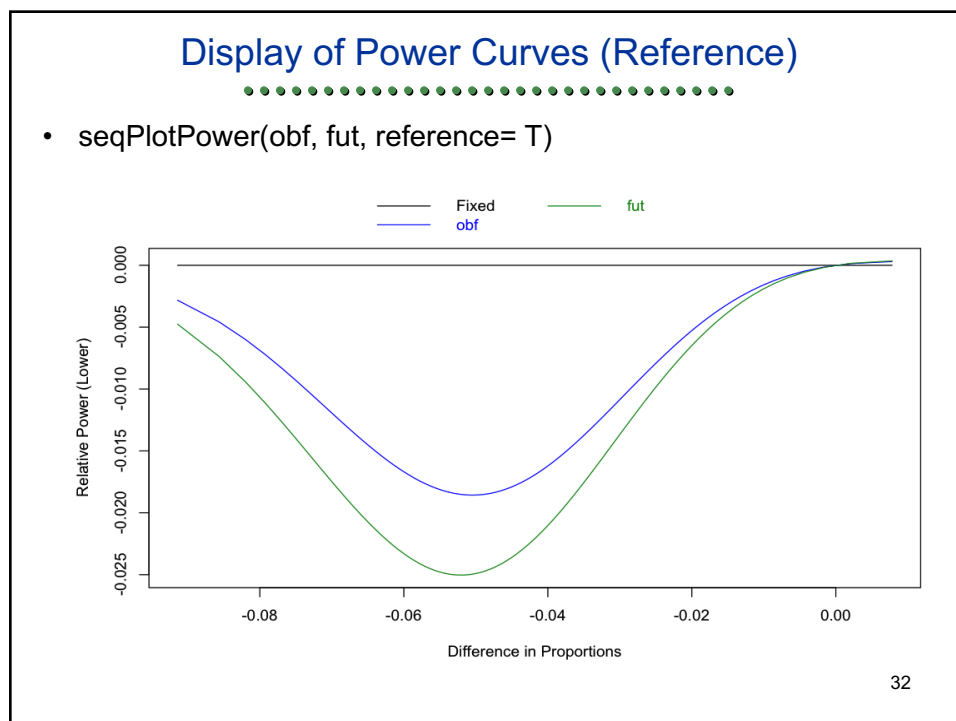


30

:



31



32

:

Plotting Sample Size Distribution

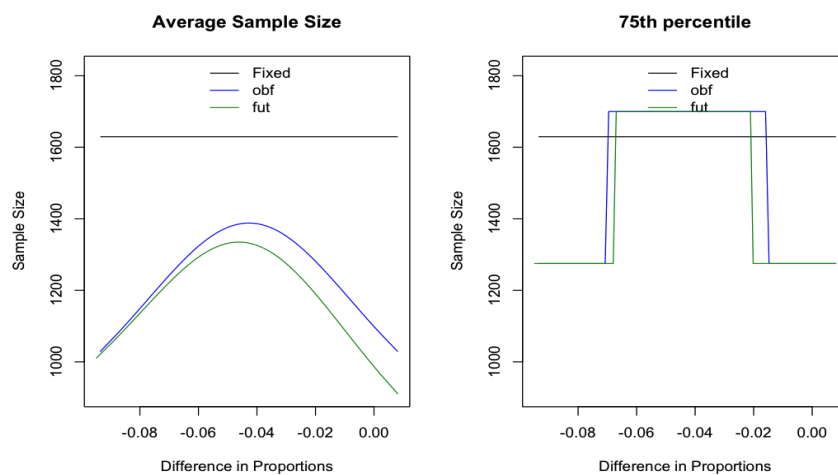
- seqPlotASN()
 - Arbitrary number of designs
 - Average sample N (ASN) and arbitrary quantiles
 - Default is 75th percentile
 - By default, include fixed design with same power as first specified design
- Axes
 - Y-axis: Statistical Information (average or quantile)
 - X-axis: Treatment effect (θ)
 - Interpretation based on probability model

33

33

Display of ASN Curves

- seqPlotASN(obf, fut)



34

34

:

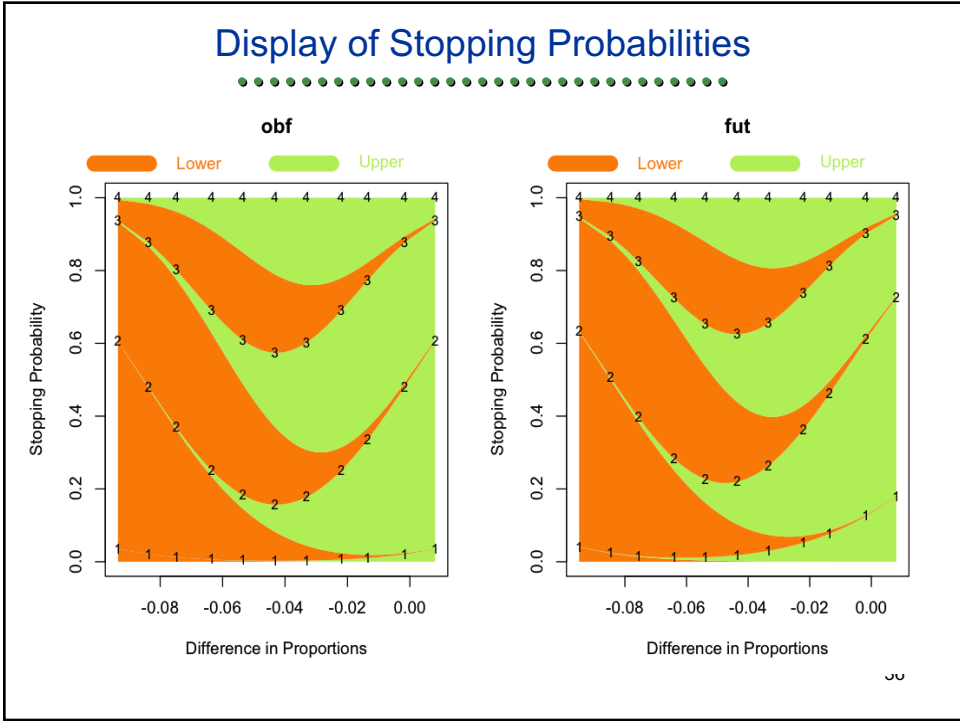
Plotting Stopping Probabilities

- seqPlotStopProb()
 - Arbitrary number of designs
 - By default, separate plots for each design
 - Different designs may have different sample sizes
 - Essentially plots of error spending functions for each alternative

- Axes
 - Y-axis: Cumulative stopping probability
 - X-axis: Treatment effect (θ)
 - Interpretation based on probability model
 - Contours labeled by analysis number
 - May not correspond to same sample size
 - Color coding by decision corresponding to stopping rule

35

35



36

Specification of Designs

Error Spending Functions

Where am I going?

The default family in RCTdesign is the Unified Family that includes a wide variety of previously described families.

RCTdesign also includes a number of designs defined on the error spending scale.

It is generally unimportant which scale is used for definition of a stopping rule, so long as they are fully evaluated

In my experience, people do not understand the scientific impact of particular error spending functions very well

37

Error Spent at Each Analysis

- seqOC() returns stopping probabilities at each analysis
- When called under the null hypothesis, this is the error spent at each analysis

```
> seqOC(fut,theta=0)
```

```
### Asymptotic Operating Characteristics  
Operating characteristics at theta= 0  
ASN= 986.6778  
Expected theta= 0.0092  
Lower Power= 0.025
```

Stopping Probabilities:

	Lower	Null	Upper	Total
Analysis time 1	0.0000	0	0.1339	0.1339
Analysis time 2	0.0024	0	0.4956	0.4981
Analysis time 3	0.0092	0	0.2713	0.2805
Analysis time 4	0.0133	0	0.0742	0.0875

38

38

:

Error Spending Functions in RCTdesign

- No matter how a design is specified, RCTdesign returns the error spending function
 - “error.spend” component of a seqDesign object has the cumulative proportion of error spent

```
> fut$error.spend
STOPPING BOUNDARIES: Error Spending Function scale
                        Efficacy Futility
Time 1 (N= 425)      0.0014  0.0341
Time 2 (N= 850)      0.0993  0.2364
Time 3 (N= 1275)     0.4684  0.5955
Time 4 (N= 1700)     1.0000  1.0000
```

- The efficacy boundary is the type 1 error spending function
- The futility boundary is the type 2 error spending function for the alternative used to define that boundary
 - (Need to fully understand the stopping boundaries)

39

39

Error Spending Families in RCTdesign

- seqDesign() argument design.family allows specification of error spending families
 - design.family="E" uses Kim & DeMets power family
 - design.family="Hwang" uses Hwang, Shih, & DeCani family
 - Includes others, as well as completely arbitrary
- Authors have described approximate correspondences to OBF and Pocock boundaries within unified family

40

40

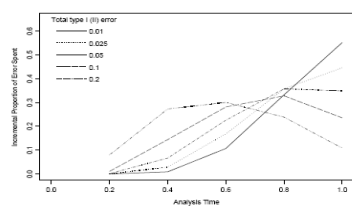
Error Spending Functions

- My view: Poorly understood even by the researchers who advocate them
- There is no such thing as THE Pocock or O'Brien-Fleming error spending function
 - Depends on type I or type II error
 - Depends on number of analyses
 - Depends on spacing of analyses

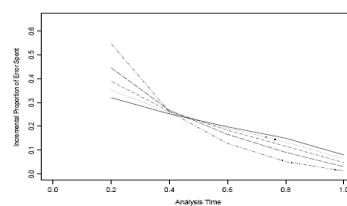
43

43

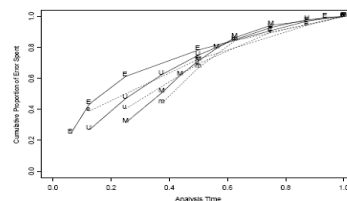
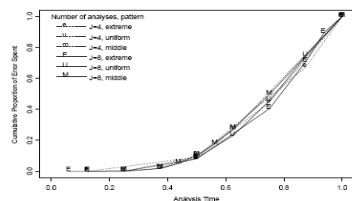
OBF, Pocock Error Spending



(a) O'Brien-Fleming Boundary Relationships



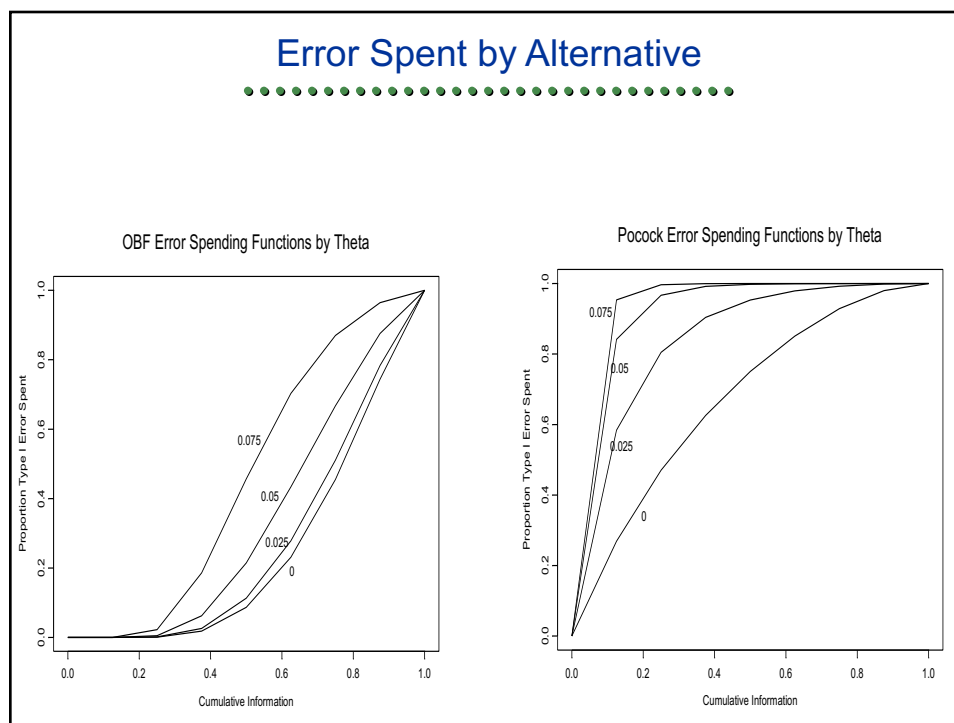
(b) Pocock Boundary Relationships



44

44

:



45

Stochastic Curtailment

.....

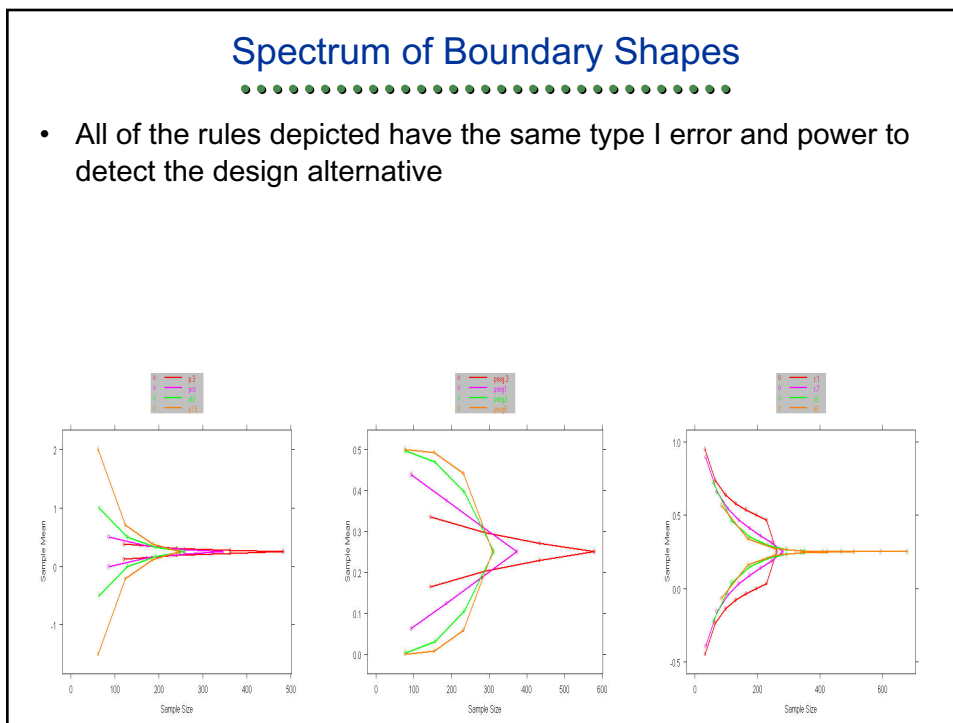
- Stopping boundaries chosen based on predicting future data
- Probability of crossing final boundary
 - Frequentist: Conditional Power
 - A Bayesian prior with all mass on a single hypothesis
 - Bayesian: Predictive Power
- Users are typically poor at guessing good thresholds
 - More on this later

46

46

Spectrum of Boundary Shapes

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



47

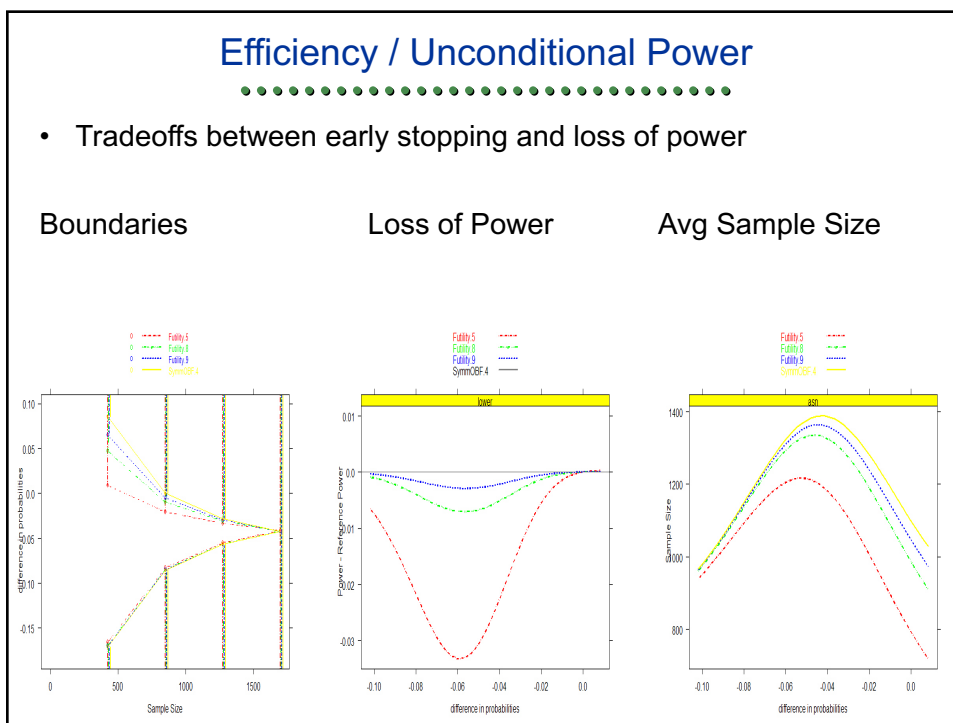
Efficiency / Unconditional Power

- Tradeoffs between early stopping and loss of power

Boundaries

Loss of Power

Avg Sample Size



48

:

Time to Event Endpoints

.....

Estimating Subject Accrual

Where am I going?

In time to event analyses, statistical information is roughly proportional to the number of events.

Additional consideration must be given to accrual of subjects.

49

Delayed Measurement of Outcome

.....

- Longitudinal studies
 - Measurement might be 6 months – 2 years after randomization
 - Interim analyses on variable lengths of follow-up
 - Use of partial data can improve efficiency (Kittelson, et al.)
- Time to event studies
 - Statistical information proportional to number of events
 - Calendar time requirements depend on number accrued and length of follow-up
- In either case: Interim analyses may occur after accrual completed
 - Group ethics of identifying beneficial treatments faster
 - Savings in calendar time costs, rather than per patient costs

50

50

:

Time to Event Endpoints

- RCTdesign allows specification of the “hazard” probability model for time to event data
 - Logrank statistic
 - Estimates of hazard ratio using Cox model
- “Sample size” computations return number of events primarily
- Additional accrual models are used to estimate
 - Number of subjects to accrue
 - Calendar time of interim analyses

51

51

Case Study: Stopping Rule

- Design of RCT to test a new drug for NSCLC
 - One-sided type 1 error: 0.025 for null of HR=1.0
 - Power: 90% to detect HR= 0.77
 - Four interim analyses with OBF efficacy, intermediate futility

```
> tte <- seqDesign("hazard", alt=0.77, power=0.9,  
+  nbr=4, P=c(1,0.8))  
> tte  
Call:  
seqDesign(prob.model = "hazard", alt.hypothesis = 0.77, nbr.analyses = 4,  
power = 0.9, P = c(1, 0.8))  
  
PROBABILITY MODEL and HYPOTHESES:  
Theta is hazard ratio (Treatment : Comparison)  
One-sided hypothesis test of a lesser alternative:  
Null hypothesis : Theta >= 1.00 (size = 0.025)  
Alternative hypothesis : Theta <= 0.77 (power = 0.900)  
  
STOPPING BOUNDARIES: Sample Mean scale  
Efficacy Futility  
Time 1 (NEv= 163.7) 0.5372 1.1891  
Time 2 (NEv= 327.4) 0.7329 0.9651  
Time 3 (NEv= 491.1) 0.8129 0.8927  
Time 4 (NEv= 654.8) 0.8561 0.8561
```

52

52

:

Case Study: Accrual (months)

- `tte <- update(tte, accrualTime= 36, studyTime= 48, eventQuantiles= 36)`

Accrual summary table:

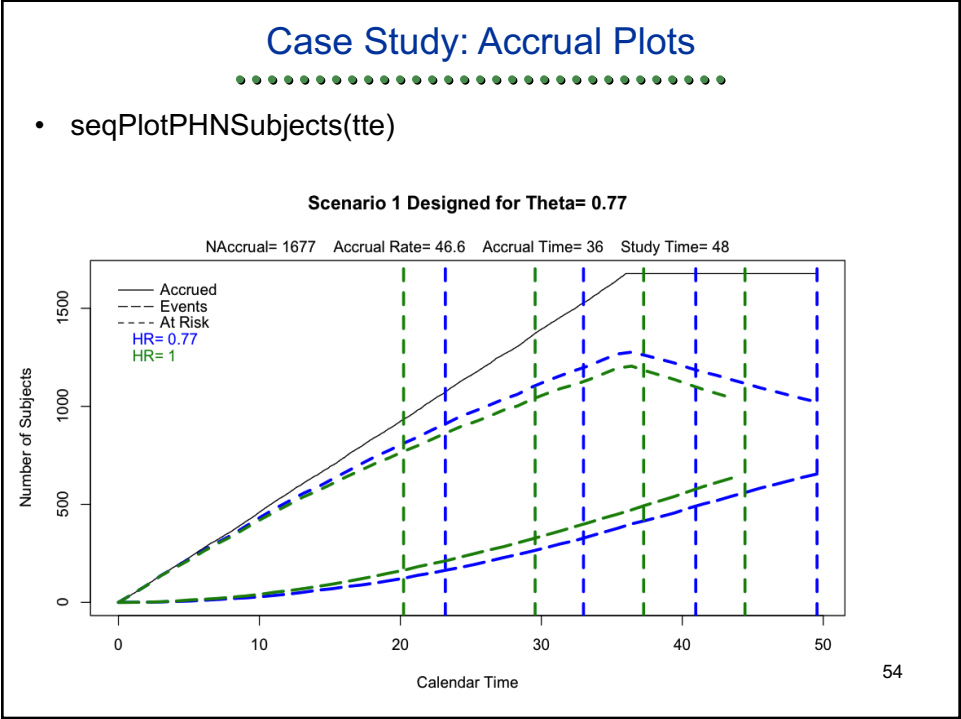
	theta	Scenario	N Accrual	accrualRate	accrualTime	studyTime
alternative	0.77	1	1677	46.58	36	48
null	1.00	1	1528	42.44	36	48

Timing of analyses:

Theta	Scenario	Analysis 1	Analysis 2	Analysis 3	Analysis 4	
0.77	1	Analysis Time	21.53	31.19	39.18	48.0
		N Accrued	1005.16	1454.43	1677.00	1677.0
		N Events	163.70	327.40	491.10	654.8
1	1	Analysis Time	21.43	31.21	39.2	48.0
		N Accrued	913.93	1324.33	1528.0	1528.0
		N Events	163.70	327.40	491.1	654.8

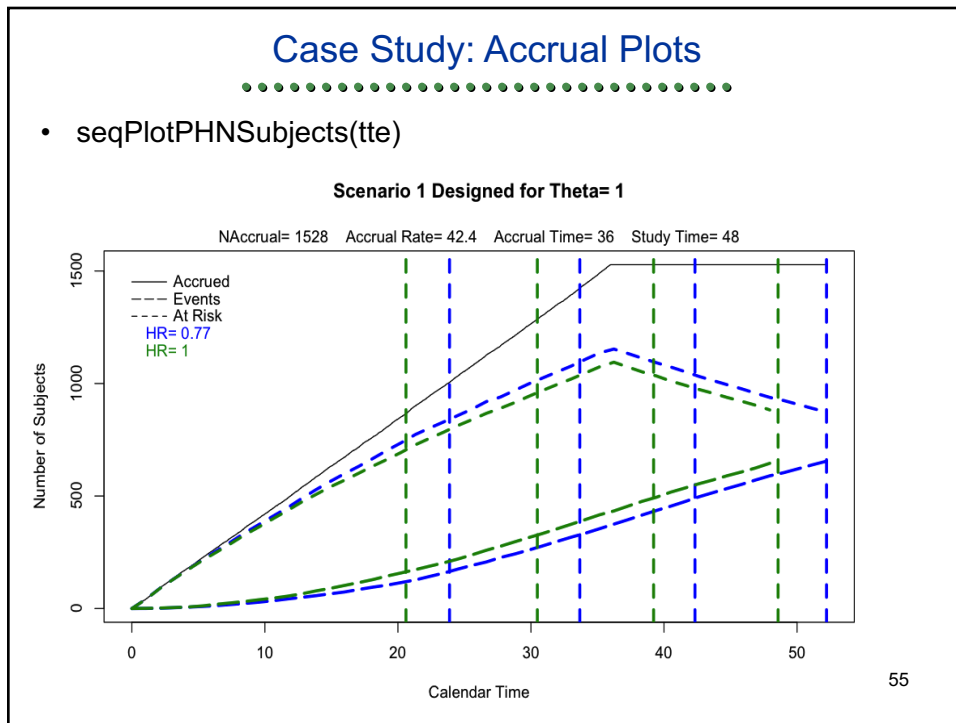
53

53



54

:



55

Implementation for Monitoring

.....

56

Flexible Determination of Boundaries

- When implementing stopping rules, must be able to accommodate changes
- Previously described methods for implementing stopping rules
 - (Adhere exactly to monitoring plan)
 - (Approximations based on design parameters: Emerson and Fleming, 1989)
 - Christmas tree approximation for triangular tests: Whitehead and Stratton, 1983
 - Error spending functions: Lan and DeMets, 1983; Pampallona, Tsiatis, and Kim, 1995
 - Constrained boundaries in unified design family: Emerson, 2000; Burington and Emerson, 2003

57

57

First Monitoring Analysis: Re-estimate N

```

> obs <- c(rbinom(200,1,0.25),rbinom(200,1,0.35))
> tx <- rep(1:0,each=200)
> mon1 <- seqMonitor(fut,obs,tx)
> mon1
Call:
seqMonitor(x = fut, response = obs, treatment = tx)

RECOMMENDATION:
  Continue

OBSERVED STATISTICS:
  Sample Size Crude Estimate Z Statistic
    400          -0.015          -0.3209

MONITORING BOUNDS:
Call:
"[for original design call, see $seqDesignCall in seqMonitor object]"

PROBABILITY MODEL and HYPOTHESES:
Theta is difference in probabilities (Treatment - Comparison)
One-sided hypothesis test of a lesser alternative:
  Null hypothesis : Theta >= 0.0000 (size = 0.025)
  Alternative hypothesis : Theta <= -0.0713 (power = 0.900)

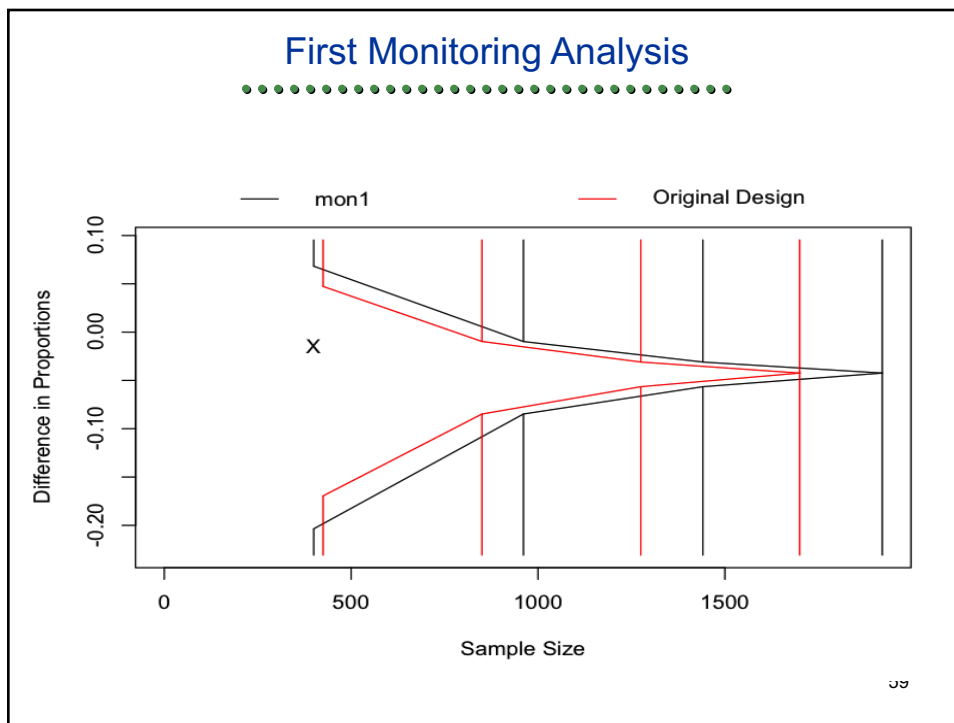
STOPPING BOUNDARIES: Sample Mean scale
              a          d
Time 1 (N= 400) -0.2036  0.0682
Time 2 (N= 961) -0.0848 -0.0098
Time 3 (N= 1441) -0.0565 -0.0310
Time 4 (N= 1921) -0.0424 -0.0424

```

58

58

:



59

Second Monitoring Analysis

```

> obs <- c(obs,c(rbinom(500,1,0.25),rbinom(500,1,0.35)))
> tx <- c(tx,rep(1:0,each=500))
> mon2 <- seqMonitor(mon1,obs,tx)
Warning message:
In seqMonitor(mon1, obs, tx) :
  1 specified future analysis time(s) were within min.increment of the current time or earlier, and are deleted
> mon2
Call:
seqMonitor(x = mon1, response = obs, treatment = tx)

RECOMMENDATION:
  Stop with decision for Lower Alternative Hypothesis

OBSERVED STATISTICS:
  Sample Size Crude Estimate Z Statistic
      400      -0.01500      -0.3209
     1400      -0.07286      -2.9544

INFERENCE:

Adjusted estimates based on observed data:
analysis.index observed      MLE      BAM      RBadj
1 -0.07285 0.001570 (-0.1212, -0.0245)

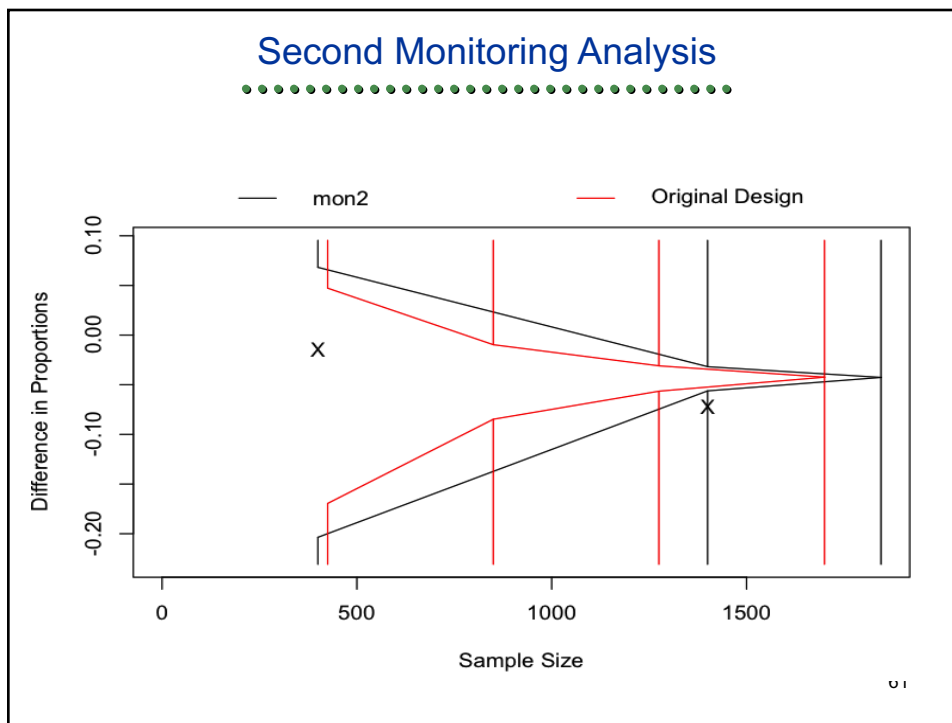
Inferences based on Analysis Time Ordering:
      MUE P-value      **** CI ****
1 -0.07285 0.001570 (-0.1212, -0.0245)

Inferences based on Mean Ordering:
      MUE P-value      **** CI ****
1 -0.07235 0.001585 (-0.1207, -0.0243)
    
```

60

60

:



61

- ### Additional Resources
-
- www.RCTdesign.org
 - SS Emerson, DL Gillen, JM Kittelson, GP Levin, SC Emerson
 - Software
 - Documentation
 - Tutorials
 - Extensions (Bayesian evaluation; adaptive design evaluation)
 - Learning
 - Short courses
 - Research talks
 - Case studies
 - Methodology
 - Technical reports on a variety of RCT-related topics

62

62