## **Statistically Speaking Lecture Series**

Sponsored by the Biostatistics Collaboration Center

Interim Analyses in Clinical Trials: An inside "Look"

Lauren Balmert Bonner, PhD

**Assistant Professor** 

Jody D. Ciolino, PhD

Associate Professor

Department of Preventive Medicine – Biostatistics

## **Biostatistics at NU**

Overview

Division of Biostatistics (Chief: Denise Scholtens), Department of Preventive Medicine (Chair: Donald Lloyd-Jones)



## **Biostatistics Centers and Cores**

#### Overview



Morthwestern Medicine<sup>®</sup> Feinberg School of Medicine

#### **Disclaimer:**

The views presented today are our own. They do not represent those of the statistical community nor those of Northwestern University at large.

We have no relevant conflicts of interest.

Inspired by a manuscript under review with Journal of Clinical and Translational Science

Coauthor: Dr. Alex Kaizer at University of Colorado

## Goals for today

- Clarify what we often mean by "interim analysis" and "interim monitoring"
- Review types of common interim analysis methods and discuss reasoning/use cases for each
  - Efficacy
  - Futility
  - Safety
  - Sample size re-estimation
- Questions and Answers / Open Discussion

## Interim Monitoring vs. Analysis

These terms are often confused + there are a lot of variations of meaning even within each overarching topic

- For purposes of our discussion...
- Interim monitoring:
  - Processes and systems important for study conduct and reporting
  - **Data integrity** and general cleanliness = primary focus
  - There is almost always an **ethical/safety component** when we talk about human subjects' research

## Interim Monitoring vs. Analysis

These terms are often confused + there are a lot of variations of meaning even within each overarching topic

• For purposes of our discussion...

- Interim **analysis**:
  - Referring to statistical tools used to guide study design modifications (most often revolving around recruitment targets) – "go/no-go", stopping bounds, adding arms, removing arms, sample size reestimation, etc.
  - Need not involve a formal hypothesis test (often does, not always!)

## Reasoning for interim monitoring + analyses

#### **Interim monitoring**

Ensuring trial conduct according to protocol, **efficiently** and ethically

- Consent verification
- Data quality checks
- Process measures screening rates, dropout rates, adherence, etc.
- **Safety** and major event reporting adverse events, deviations, etc.

Findings on interim monitoring help drive study conduct decisions.

#### **Interim analysis**

Ensuring the study is set up for success (it addresses the research question + is equipped to address the research question) with an adequate risk-benefit profile

- Making **efficient** use of participant time, data, and general resources
- Continual assessment of safety signals

Findings on interim analysis help drive study design considerations.

Title: Guidance on interim analysis methods in clinical trials

Authors: Jody D. Ciolino, PhD<sup>1</sup>; Alexander M. Kaizer, PhD<sup>2</sup>; Lauren Balmert Bonner, PhD<sup>1</sup>

#### Affiliations:

<sup>1</sup> Department of Preventive Medicine (Biostatistics), Northwestern University Feinberg School of Medicine, Chicago, Illinois

<sup>2</sup> Department of Biostatistics & Informatics, Colorado School of Public Health, Aurora, Colorado

- The term "interim analysis" in clinical trials has multiple meanings.
- In general, interim analyses help **guide decisions on overall clinical trial** modifications, specifically those pertaining to the study sample size or recruitment targets
- We often encounter a lot of confusion and misunderstanding when it comes to interim analysis

#### Table 1. Summary of Interim Analysis Types

	Explanation	Justification for Use
Efficacy	<ul> <li>Early termination of a trial that is showing promising results</li> <li>Control of type I error through group sequential methods or alpha spending functions</li> </ul>	<ul> <li>Usually for longer, larger studies and later phases of research</li> <li>Ethical imperative for a promising treatment to reach the entire target clinical population</li> </ul>
Futility	<ul> <li>Early termination of a trial that is not likely to achieve the intended objective (e.g., little chance of finding a "significant" treatment effect at the end of the study)</li> <li>Employed through group sequential methods, error spending functions, conditional power, or predictive power</li> </ul>	<ul> <li>Reduces costs, resources, and patient burden for a trial with a low probability of "success"</li> <li>Usually for mid-late phase studies</li> <li>Helpful in the context of recruitment and retention challenges</li> </ul>
Safety	<ul> <li>Early termination (or pausing) of a trial for safety concerns</li> <li>Should be coupled with efficacy analyses to evaluate the benefit-to-risk ratio</li> </ul>	<ul> <li>Incorporated across all phases of research</li> <li>Particularly important for vulnerable populations and high-risk interventions with more "serious" outcomes (e.g., death)</li> </ul>
Sample size re-estimation	<ul> <li>Reassessment of the sample size required to ensure adequate power using updated information from interim trial data</li> <li>Can be blinded or unblinded</li> <li>May not necessarily spend alpha</li> </ul>	<ul> <li>Allows for interim look at assumptions (standard deviations, event rates, correlations, etc.)</li> <li>May be particularly useful for midlate phase studies</li> </ul>

#### Morthwestern Medicine®

Feinberg School of Medicine

Efficacy

- Involves a statistical hypothesis test evaluating one arm over another
- If there is a "large enough" signal early on in the study, it may be ethically imperative and most efficient to stop the study early
- Threshold for sufficient evidence is subject to debate and topic of statistical research and literature
- We cannot simply use a "statistically significant" finding to guide this decision...

- Recall: type I error = probability of finding a significant result [usually p<0.05] when in fact we should not, as there is no effect
- The more times we "look" at the data (i.e., conduct statistical tests), the more likely we are to find a significant result (i.e., make a type I error)



Type I error illustration for hypothetical null effect trial

- Consider a hypothetical, two-arm clinical trial
- Binary outcome = "success" of intervention
- Plan to conduct a statistical test at the 5% level of significance
- We can simulate no underlying difference between the two study arms
  - Assume p(success) = 0.30 across arms
  - N=355 planned per arm
  - Simulate under null (H0)
  - IF we conduct a statistical test and find p<0.05, we are making a type I error
- There are countless ways the test statistic could "behave" over the course the study, but consider one hypothetical scenario...

Type I error illustration for hypothetical null effect trial



Progression through Trial, Increasing Information

% of Participants Enrolled with Outcome Data Observed	20%	40%	41%**	60%	80%	100%
N per Arm	71	142	143	213	284	355
p-value	0.019	0.057	0.042	0.185	0.096	0.212
Arm 1 Successes	11	29	29	49	73	93
Arm 2 Successes	24	44	45	62	92	109
Difference in Successes (N)	13	15	16	13	19	16
Difference in Proportions	0.183	0.105	0.112	0.061	0.067	0.045

#### \*\*between 40-41%, included for illustrative purposes

Illustration of test statistic behavior with increasingly large sample sizes



Increasing Information, Thousands of Participants per Arm

- Recall: type I error = probability of finding a significant result [usually p<0.05] when in fact we should not as there is no effect
- The more times we "look" at the data (i.e., conduct statistical tests), the more likely we are to find a significant result (type I error)
- Utilize **methodology to control type I error** with repeated looks at the data
- These methods help us decide whether the results at an interim analysis are "significant enough" to warrant early stopping



Common methods for controlling type I error

- Another two-arm hypothetical trial
- 4 interim looks + 1 final analysis
- Typical "group sequential" stopping bounds look like this...



Common methods for controlling type I error

- **Pocock** bounds have a constant threshold (p<0.016) for all 5 analyses
- **Peto** bounds have a stringent (p<0.001) threshold for the first 4, then uses p<0.05 at the end
- **O'Brien-Fleming** (perhaps most common) stringent thresholds early on, and increasingly less stringent as time goes on; final analysis uses p<0.04



- Note: "alpha spending functions" incorporate these ideas but allow for more flexibility
- Can accommodate different ways of "spending" type I error (differing weights), and the timing of analyses needed not be evenly spaced
- Sometimes these terms "group sequential" and "alpha spending" are used interchangeably since they are so closely linked

# Thrombectomy for Stroke in the Public Health Care System of Brazil

Martins SO, Mont'Alverne F, Rebello LC, et al. Thrombectomy for stroke in the public health care system of Brazil. *New England Journal of Medicine*. 2020;382(24):2316-2326

- Randomized stroke patients to standard-of-care (SOC) or SOC + mechanical thrombectomy
- Proposed interim analysis using information from 90-day follow-up
- Primary outcome = modified Rankin scale (measure of disability) at 90 days



# Thrombectomy for Stroke in the Public Health Care System of Brazil

Martins SO, Mont'Alverne F, Rebello LC, et al. Thrombectomy for stroke in the public health care system of Brazil. *New England Journal of Medicine*. 2020;382(24):2316-2326

- Results at first interim analysis (N=174):
  - OR = 2.24 (1.30, 3.88) p-value = 0.004, in favor of thrombectomy
- Data and Safety Monitoring Board (DSMB) recommended early stopping
- At time of early termination, N=221 had been randomized and were included in final analyses
  - OR = 2.28 (1.41, 3.69) p-value = 0.001



#### Interim Analysis for Efficacy: Implications

- Controlling type I error in any trial is important
- Things get more complicated if we plan to incorporate interim analyses for efficacy
- Group sequential methods and alpha spending functions allow researchers a tool to maintain control over type I error
- In the example trial, investigators were able to address their research question with fewer participants than planned → more efficient use of participant time and study resources
  - This is the heart of the reasoning behind these analyses
  - Caveat intervention being studied should be in later stages (e.g., phase III clinical trial) of development; must consider the big picture



#### Morthwestern Medicine®

Feinberg School of Medicine

Futility

iody.ciolino@northwestern.e28

#### **Futility Analyses**

- Stopping a trial for futility suggests that observing a statistically significant result at the end of the study is unlikely
- This can **increase efficiencies** with respect to cost, resources, and participant burden
- There are similar methods to the group sequential methods for futility (or error spending functions)



#### Figure 1 A futile action can never achieve its goals.

Kite S, Wilkinson S. Beyond futility: to what extent is the concept of futility useful in clinical decision-making about CPR? Lancet Oncol. 2002 Oct;3(10):638-42. doi: 10.1016/s1470-2045(02)00878-1. PMID: 12372726.



One-sided test with efficacy and futility stopping bounds



# The Stroke and Hyperglycemia Insulin Network Effort (SHINE)

Johnston KC, Bruno A, Pauls Q, et al. Intensive vs Standard Treatment of Hyperglycemia and Functional Outcome in Patients With Acute Ischemic Stroke: The SHINE Randomized Clinical Trial. JAMA. 2019;322(4):326-335. doi:10.1001/jama.2019.9346

- Randomized trial to evaluate efficacy of intensive glucose control during ischemic stroke (N=1400)
- Non-binding futility thresholds using error spending function approach
- 4<sup>th</sup> interim analysis (N=1151), trial was stopped for futility
- Final results: no significant difference in proportion with 90-day favorable outcome between groups (20.5% vs. 21.6%)

Interim Analysis Sample Size	P-value Futility Threshold
500	0.949
700	0.896
900	0.652
1100	0.293





- Probably a bit more common are **conditional power or predictive power** approaches...
- In general: power = probability(reject null | some assumption with respect to true underlying effect)

		True	
		Null is True (No Difference)	Null is not True (Difference)
Tost	Reject Null	Type I Error	Power
iest	Fail to Reject Null	Confidence	Type II Error

 Conditional power = probability(reject null | data observed up to this point in the trial and assumption with respect to effect)

#### **Conditional Power**

- Based on...
  - Information fraction (how far we are through the trial)
  - Current test statistic from interim analysis
  - Assumption of effect at end of trial
    - Current trend
    - H1 (from beginning of trial)
    - H0 or null effect
  - If this calculated conditional power is low (usually below 20%, 15%, 10%), then we might consider stopping a study

#### **Predictive Power**

- Bayesian alternative to conditional power
- Likelihood of demonstrating treatment efficacy at the end of the study
- Estimated by...
  - Updating prior assumptions with observed data
  - Average conditional power over this distribution
- Avoids having to assume a specific treatment effect (as in the conditional power approach)

## Vesnarinone in Heart Failure Trial (VEST)

Cohn JN, Goldstein SO, Greenberg BH, et al. A dose-dependent increase in mortality with vesnarinone among patients with severe heart failure. N Engl J Med 1998;339:1810–1816.

- Patients with heart failure randomized to placebo or Vesnarinone
- Primary outcome: All-cause mortality
- Sample size calculations:
  - 20% 1-year mortality in placebo
  - Power to detect 30% reduction (0.2-0.2\*0.3 = 0.14)
  - Sample size of 3618 would provide 90% power

## Vesnarinone in Heart Failure Trial (VEST)

Cohn JN, Goldstein SO, Greenberg BH, et al. A dose-dependent increase in mortality with vesnarinone among patients with severe heart failure. N Engl J Med 1998;339:1810–1816.

- Patients with heart failure randomized to placebo or Vesnarinone
- Primary outcome: All-cause mortality
- Sample size calculations:
  - 20% 1-year mortality in placebo
  - Power to detect 30% reduction (0.2-0.2\*0.3 = 0.14)
  - Sample size of 3618 would provide 90% power

But the trial was not stopped early – due to external trial information

<b>Fable 17.2</b> Conditional		Informatio	Information fraction		
Heart Failure Trial (VEST)	RR	0.50	0.67	0.84	
[129]	0.50	0.46	< 0.01	< 0.01	
	0.70	0.03	< 0.01	< 0.01	
	1.0	< 0.01	< 0.01	< 0.01	
Considered a range of effects:	1.3	< 0.01	< 0.01	< 0.01	
-Beneficial effects: 0.5-0.7	1.5	< 0.01	< 0.01	< 0.01	
-Null effect: 1.0 -Negative trends: 1.3-1.5	RR = rela	utive risk			
M Northern stown Madicine*		Fried	dman LM. Furberg CD. DeMets	DL. Reboussin D. Gra	

Feinberg School of Medicine

Friedman LM, Furberg CD, DeMets DL, Reboussin D, Granger GB. *Fundamentals of Clinical Trials*. Fifth Edition. 2015.

## **Futility Analysis: Implications**

- Incorporating a futility assessment can **increase efficiency of the trial**, allowing trials that are unlikely to meet their objectives to stop early ultimately reducing costs, preserving resources, and limiting patient burden
- Particularly in large clinical trials or vulnerable patient populations, an interim futility assessment may be essential to prevent patients from being unnecessarily randomized to ineffective treatments
- Despite common misconceptions, an interim analysis incorporating a futility assessment **alone does not inflate the type I error**
- It can, however, have **implications on type II error**, reducing the overall power of the study by stopping early



- Potential to pose challenges in interpretations smaller than expected sample size decreases precision around treatment effect estimates
- Subgroup analyses and analyses examining heterogeneity of intervention effects will inevitably suffer (they are typically underpowered already)



#### Morthwestern Medicine®

Feinberg School of Medicine

Safety

iody.ciolino@northwestern.e84

## Interim analysis for safety



- Note: all clinical studies (regardless of phase) should incorporate safety monitoring
- It is impossible to foresee all potential safety issues ahead of time, and **safety analyses are seldom adequately powered**
- Unanticipated safety issues should merit consideration for early stopping; for example,
  - Serious and unanticipated events or event rates
  - Events that are related and unexpected



## Interim analysis for safety

 Discussion on formal interim safety analyses should consider the benefit-to-risk assessment: i.e., it should be paired with an interim efficacy analysis

• **Context is key** in this assessment.... For example, a cardiovascular secondary prevention trial to prevent subsequent myocardial infarction may have the expectation of some myocardial infarction events, whereas it may be extremely concerning to observe the same events in a behavioral intervention in a generally healthy population.



## The EARLY Trial

Sperling R, Henley D, Aisen PS, et al. Findings of efficacy, safety, and biomarker outcomes of atabecestat in preclinical Alzheimer disease: a truncated randomized phase 2b/3 clinical trial. JAMA neurology. 2021;78(3):293-301

- Randomized phase 2b/3 trial assessing effects of atabecestat in preclinical Alzheimer's disease
  - 3 arms (2 doses of active intervention + placebo)
  - Double-blind
  - Primary outcome: change from baseline in cognitive composite score
- Plan:
  - N=1650 across 143 sites
  - Interim analysis for futility
  - NO formal interim efficacy analyses





•38

#### Interim analysis for safety: Implications

- Monitoring safety outcomes and adverse events is crucial for maintaining study integrity and protecting study participants
- Decision to stop a trial early for safety concerns should generally be made in the context of the **benefit-to-risk ratio** 
  - No "one size fits all" approach to assessing benefit-to-risk ratio

Remember: **Context is key** in this assessment.... For example, a cardiovascular secondary prevention trial to prevent subsequent myocardial infarction may have the expectation of some myocardial infarction events, whereas it may be extremely concerning to observe the same events in a behavioral intervention in a generally healthy population.

#### Interim analysis for safety: Implications

- Safety monitoring may raise concerns about multiple "looks" at the data
  - When coupled with efficacy analysis, incorporate conservative alpha-spending approach
- Safety concerns may warrant a temporary pause to further assess causality – consider whether to pause vs. completely terminate the study



#### Morthwestern Medicine®

Feinberg School of Medicine

#### Sample Size Re-estimation

jody.ciolino@northwestern.edu

- Designed to modify the planned sample size based on the accumulating data
- Accounts for uncertainty when conducting power calculations during the initial planning of the study



Morthwestern Medicine® Feinberg School of Medicine

- Designed to modify the planned sample size based on the accumulating data
- Accounts for uncertainty when conducting power calculations during the initial planning of the study



Morthwestern Medicine\* Feinberg School of Medicine

- Approaches facilitate a revised sample size calculation using information for the ongoing assessment of event rates, the estimation of nuisance parameters (e.g., the variance of a continuous outcome), or the effect size expected
- Re-estimating a sample size at an interim stage can increase the likelihood of a successful trial, but may result in a substantial increase in the needed sample size if the initial sample size assumptions were very different from what is observed



Feinberg School of Medicine





- Two categories: blinded or unblinded
- Blinded: used to revise estimation of nuisance parameters (e.g., variance)
   often used pooled estimates
- **Unblinded**: based on comparative interim results; ideal when uncertainty about estimates of effect size and nuisance parameters allows for capturing an effect that may still be clinically meaningful, but differs from original assumption.

#### Tenecteplase versus Alteplase before Endovascular Therapy for Ischemic Stroke (EXTEND-IA TNK) Trial

Campbell BC, Mitchell PJ, Churilov L, et al. Tenecteplase versus alteplase before thrombectomy for ischemic stroke. New England Journal of Medicine. 2018;378(17):1573-1582.

- Non-inferiority randomized (1:1) trial enrolling patients with ischemic stroke eligible to undergo intravenous thrombolysis and endovascular thrombectomy
- Primary outcome: Proportion of participants with restoration of blood flow to >50% of the affected arterial territory
- Initial sample size calculation: **120 participants to provide 80% power**
- Proposed blinded sample size re-estimation at n=100 using conditional power approach

#### Tenecteplase versus Alteplase before Endovascular Therapy for Ischemic Stroke (EXTEND-IA TNK) Trial

Campbell BC, Mitchell PJ, Churilov L, et al. Tenecteplase versus alteplase before thrombectomy for ischemic stroke. New England Journal of Medicine. 2018;378(17):1573-1582.

- Non-inferiority randomized (1:1) trial enrolling patients with ischemic stroke eligible to undergo intravenous thrombolysis and endovascular thrombectomy
- Primary outcome: Proportion of participants with restoration of blood flow to >50% of the affected arterial territory
- Initial sample size calculation: **120 participants to provide 80% power**
- Proposed blinded sample size re-estimation at n=100 using conditional power approach
- Conditional power: 0.52 -> Increased sample size to N=202 to maintain adequate power

## Sample size re-estimation: Implications

• Goal to prevent underpowered studies



- Be careful as this **may have impact on type I error** without using appropriate methods to control
- What if re-estimated sample size is not feasible?
- What if re-estimated sample size is based on an observed effect that is **no longer clinically meaningful**?
- Care should be taken in **how results from interim re-estimation are reported** for ongoing studies -> may be possible to back-calculate the effect size!

#### M Northwestern Medicine®

Feinberg School of Medicine

#### Tying it all together

iody.ciolino@northwestern.e52

#### Table 1. Summary of Interim Analysis Types

	Explanation	Justification for Use
Efficacy	<ul> <li>Early termination of a trial that is showing promising results</li> <li>Control of type I error through group sequential methods or alpha spending functions</li> </ul>	<ul> <li>Usually for longer, larger studies and later phases of research</li> <li>Ethical imperative for a promising treatment to reach the entire target clinical population</li> </ul>
Futility	<ul> <li>Early termination of a trial that is not likely to achieve the intended objective (e.g., little chance of finding a "significant" treatment effect at the end of the study)</li> <li>Employed through group sequential methods, error spending functions, conditional power, or predictive power</li> </ul>	<ul> <li>Reduces costs, resources, and patient burden for a trial with a low probability of "success"</li> <li>Usually for mid-late phase studies</li> <li>Helpful in the context of recruitment and retention challenges</li> </ul>
Safety	<ul> <li>Early termination (or pausing) of a trial for safety concerns</li> <li>Should be coupled with efficacy analyses to evaluate the benefit-to-risk ratio</li> </ul>	<ul> <li>Incorporated across all phases of research</li> <li>Particularly important for vulnerable populations and high-risk interventions with more "serious" outcomes (e.g., death)</li> </ul>
Sample size re-estimation	<ul> <li>Reassessment of the sample size required to ensure adequate power using updated information from interim trial data</li> <li>Can be blinded or unblinded</li> <li>May not necessarily spend alpha</li> </ul>	<ul> <li>Allows for interim look at assumptions (standard deviations, event rates, correlations, etc.)</li> <li>May be particularly useful for midlate phase studies</li> </ul>

## **Considerations for Interim Analyses**

- Pre-specify as much as possible
- Describe anticipated timing, proposed methodology, pre-specified rules to guide decisions
- Timing of analysis can be flexible
- Often specified when some proportion of participants is enrolled and meet a particular study milestone (e.g., 50% of participants completed 6-week follow-up)
- Balance between maximal information (later interim analysis) vs. ensuring adequate time to make any modifications and reducing potential risk to participants as much as possible

## Considerations for interim analysis

- Think about potential logistical implications...if an interim sample size re-estimation is proposed, are there adequate resources to support an increase in sample size if indicated?
- Evaluation of interim analysis results **should not be interpreted in isolation**, but rather in the context of other internal study factors and external contemporaneous issues
- Any interim analysis results and statistical tools are intended to serve as **guidelines**
- <u>Transparency</u> in disseminating trial results when interim analyses were conducted is also critical



- Interim monitoring of data quality and integrity ≠ interim analysis to guide study design modifications
- Both **require coordination and pre-specification** of protocol and procedural elements to the extent possible

- Monitoring does not have implications for type I error issues; interim analyses can, but it is not always the case
- Remember that there is no "one size fits all" + recommendations should be made with a big picture view

## Relying on a Data and Safety Monitoring Board (DSMB)

A minor tangent that is important to consider in interim analyses

- A DSMB is an independent group of experts (and potentially patient representative[s]), appointed by a sponsor/investigator to review accumulating <u>clinical trial</u> data on a regular basis.
- The external DSMB will be instrumental in reviewing these data and any analysis results to <u>guide recommendations</u> (remember they are recommendations)
- Interim results should be kept strictly with the DSMB and the unblinded study statistician. Only high-level recommendations from the DSMB and/or modifications to the trial should be communicated to the study team or external entities while the study is ongoing.

#### **DSMB** Composition

- Usually 3-5 total members, can be more
  - Clinical expert in relevant field(s)
  - Biostatistician
  - Basic scientist
  - Regulatory specialist
  - Ethicist
  - Patient representative

#### **DSMB Review Process**



#### DSMB Review Process – Meeting Format



Use caution when discussing decisions based on interim analyses

Morthwestern Medicine\*

Thank you for your attention today + please feel free to reach out with comments/questions!

## Upcoming Statistically Speaking Lectures

Thursday, April 6 12-1pm Longitudinal Data: Challenges and Opportunities

Nathan Gill, PhD, Assistant Professor, Division of Biostatistics, Department of Preventive Medicine *Location: Lurie - Hughes Auditorium & Zoom\** 

## Thank you!



#### Expertise in biostatistics, statistical programming and data management

Since 2004, the Biostatistics Collaboration Center (BCC) has partnered with Northwestern investigators at every level – from residents and postdoctoral fellows to junior faculty and well-established senior investigators.



#### Meet Our Team

Our faculty and staff have collaborated with basic science, clinical, and health services investigators in over 60 Northwestern units across Chicago and Evanston campuses, and with all of NU's clinical partners. Meet our team and learn more about our expertise.

