

Group Sequential Designs with Early Stopping for Efficacy and Futility

Christopher Jennison

Department of Mathematical Sciences,
University of Bath, UK

<http://people.bath.ac.uk/mascj>

**NHLBI Workshop on Clinical Trial Designs:
Innovative Endpoints and Futility Monitoring**

Bethesda

September 2023

Group Sequential Designs

Group sequential designs allow a clinical trial to be terminated

For efficacy:

When there is overwhelming evidence that the new treatment is effective,

For futility:

When it is clear the trial is unlikely to reach a positive conclusion.

In retrospective analyses of 72 ECOG cancer studies, Rosner & Tsiatis (*Statist. in Med.*, 1989) found that, had group sequential stopping rules been applied, early stopping (mostly to accept H_0) would have occurred in around 80% of cases.

Many clinical trials have a formal stopping rule for efficacy but informal guidelines for futility stopping. Why should these issues be treated differently — and is this a good idea?

Outline of talk

1. Defining an efficacy stopping rule through an error spending function.
2. Using conditional power and predictive power to guide stopping for futility.
3. Error spending designs with an efficacy boundary and a non-binding futility boundary.

1. Defining a stopping rule for efficacy

Consider a Phase 3 clinical trial comparing a new treatment against a standard.

Let θ denote the “effect size”, a measure of the improvement in the new treatment over the standard.

We shall test the null hypothesis $H_0: \theta \leq 0$ against $\theta > 0$.

Rejecting H_0 allows us to conclude the new treatment is superior.

We allow type I error probability α for rejecting H_0 when it is true.

We specify power $1 - \beta$ as the probability that H_0 should be rejected when $\theta = \delta$.

Here δ is, typically, the minimal clinically significant treatment difference.

Joint distribution of parameter estimates

Reference: Ch. 11 of *Group Sequential Methods with Applications to Clinical Trials*, Jennison & Turnbull, 2000.

Let $\hat{\theta}_k$ denote the estimate of θ based on data at analysis k .

The information for θ at analysis k is

$$\mathcal{I}_k = \{\text{Var}(\hat{\theta}_k)\}^{-1}, \quad k = 1, \dots, K.$$

Canonical joint distribution of $\hat{\theta}_1, \dots, \hat{\theta}_K$

In many situations, $\hat{\theta}_1, \dots, \hat{\theta}_K$ are approximately multivariate normal,

$$\hat{\theta}_k \sim N(\theta, \mathcal{I}_k^{-1}), \quad k = 1, \dots, K,$$

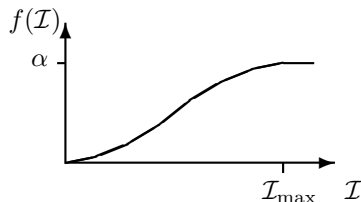
and

$$\text{Cov}(\hat{\theta}_{k_1}, \hat{\theta}_{k_2}) = \text{Var}(\hat{\theta}_{k_2}) = \mathcal{I}_{k_2}^{-1} \quad \text{for } k_1 < k_2.$$

Error spending tests (Jennison & Turnbull, Ch. 7)

When the sequence $\mathcal{I}_1, \mathcal{I}_2, \dots$ is unpredictable, a group sequential design must adapt to observed information levels. Lan & DeMets (*Biometrika*, 1983) introduced “error spending” tests of $H_0: \theta = 0$ against $\theta \neq 0$.

Maximum information design with error spending function $f(\mathcal{I})$



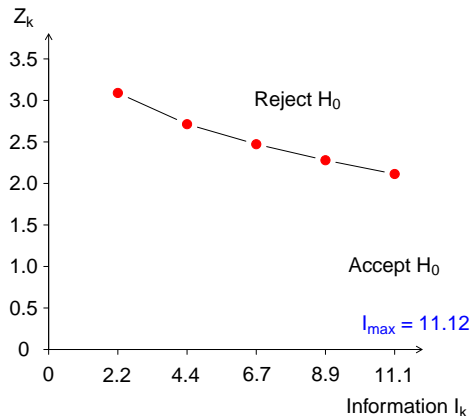
The boundary at analysis k is set to give cumulative type I error probability (under $\theta = 0$) equal to $f(\mathcal{I}_k)$.

If the target information level, \mathcal{I}_{\max} , is reached without rejection of H_0 , then H_0 is accepted.

Example: A trial design with an efficacy boundary only

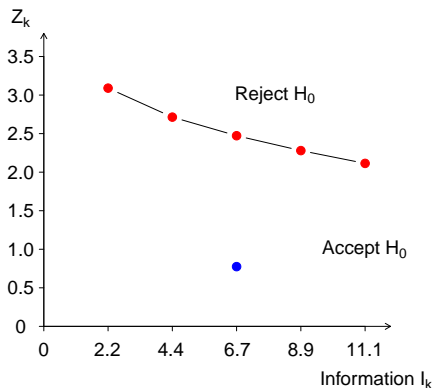
A test with 5 planned analyses, type I error probability $\alpha = 0.025$, power 0.9 if $\theta = \delta = 1$, and type I error spending function

$$f(\mathcal{I}) = \min\{(\mathcal{I}/\mathcal{I}_{\max})^2, 1\} \alpha.$$



2. Using conditional power to guide stopping for futility

Suppose the trial has reached analysis 3, $\hat{\theta}_3 = 0.3$ and $Z_3 = 0.78$.



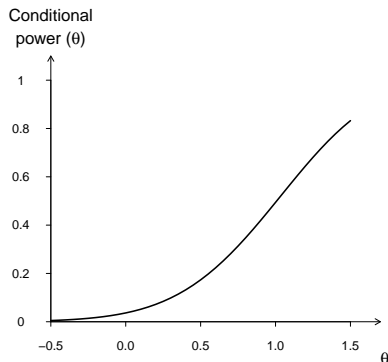
One may ask

“How likely is it that the trial’s final outcome will be positive?”

Using conditional power to guide stopping for futility

We can compute the “conditional power function”,

$$P_{\theta}\{H_0 \text{ will be rejected at analysis 4 or 5} \mid Z_3 = 0.78\}.$$

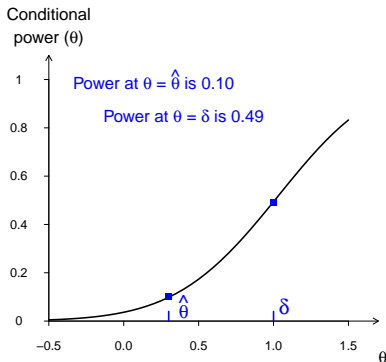


However, we do not know the true value of θ .

Using conditional power to guide stopping for futility

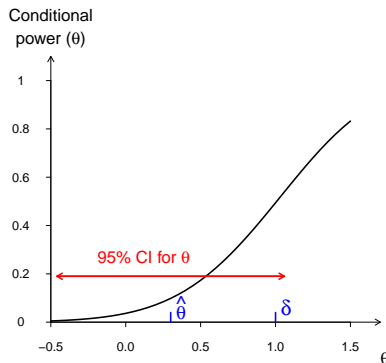
One may focus on the conditional power assuming θ is equal to the current estimate, $\hat{\theta}_3 = 0.3$.

Or, one may focus on conditional power assuming $\theta = \delta = 1$, the value used in the power calculation.



Using conditional power to guide stopping for futility

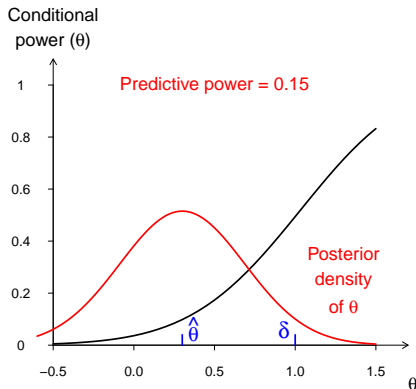
It is important to remember that the estimate of θ at an interim analysis has a high variance.



Adopting a Bayesian approach, one can integrate conditional power over a posterior distribution to obtain a “predictive power”.

Using predictive power to guide stopping for futility

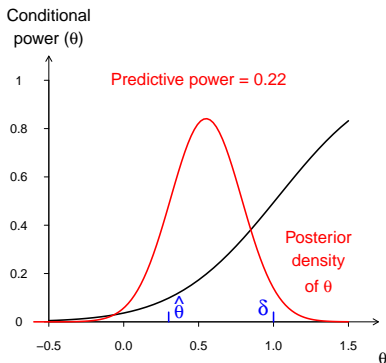
A common practice is to assume a flat (improper) prior for θ in calculating predictive power.



In this case the posterior distribution of θ is $N(0.3, 0.39^2)$.

Using predictive power to guide stopping for futility

Given the high variance of the interim estimate $\hat{\theta}_3$, the choice of prior can have a significant impact on the predictive power.



Under the prior $\theta \sim N(0.7, 0.3^2)$, the posterior distribution of θ is $\theta \mid \hat{\theta}_3 \sim N(0.55, 0.24^2)$, and predictive power rises to 0.22.

Using predictive power to guide stopping for futility

Once you have calculated your chosen conditional power or predictive power, the question remains:

How high should the conditional probability of success be to justify continuation of the trial?

Here, one needs to balance

The benefits from saving resources in this study and moving on to conduct trials for other promising therapies,

The risk of stopping the current trial prematurely when it would have gone on to produce a positive result.

Decision making is hard when conditional power is low but there is still a non-negligible chance the trial may still succeed.

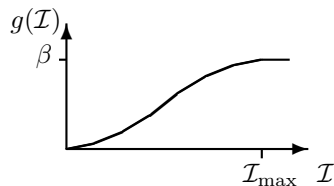
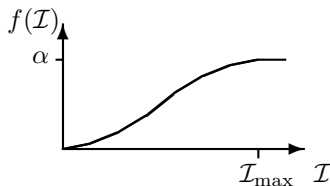
3. Error spending tests with efficacy and futility boundaries

For a one-sided test of $H_0: \theta \leq 0$ against $\theta > 0$ with

Type I error probability α at $\theta = 0$,

Type II error probability β at $\theta = \delta$,

we need two error spending functions.



Type I error probability α is spent according to the function $f(\mathcal{I})$, and type II error probability β (under $\theta = \delta$) according to $g(\mathcal{I})$.

Treating the futility boundary as “non-binding”, we calculate Type I error probabilities ignoring the futility boundary.

Error spending tests with efficacy and futility boundaries

Recall, we want a group sequential test of $H_0: \theta \leq 0$ vs $\theta > 0$ with

$$P_{\theta=0}\{\text{Reject } H_0\} = \alpha,$$

$$P_{\theta=\delta}\{\text{Accept } H_0\} = \beta,$$

Analyses at $\mathcal{I}_k = (k/K) \mathcal{I}_{\max}$, $k = 1, \dots, K$.

If we specify α , β , δ , K and \mathcal{I}_{\max} , we can find the stopping rule that minimises a criterion such as

$$\sum_i w_i E_{\theta_i}(\mathcal{I}) \quad \text{or} \quad \int w(\theta) E_{\theta}(\mathcal{I}) d\theta.$$

See:

Barber & Jennison (*Biometrika*, 2002),

Öhrn (*PhD thesis, University of Bath*, 2011),

Jennison & Turnbull (*Kuwait J. Science*, 2013).

Error spending tests with efficacy and futility boundaries

Barber & Jennison (2002) and Öhrn (2011) observe that group sequential tests with error spending functions of the form

$$f(\mathcal{I}) = \min\{(\mathcal{I}/\mathcal{I}_{\max})^{\rho_1}, 1\} \alpha \quad (\text{type I error})$$

and

$$g(\mathcal{I}) = \min\{(\mathcal{I}/\mathcal{I}_{\max})^{\rho_2}, 1\} \beta \quad (\text{type II error})$$

have high efficiency for a variety of optimality criteria.

The values of ρ_1 and ρ_2 determine \mathcal{I}_{\max} and, hence, the trial's maximum sample size.

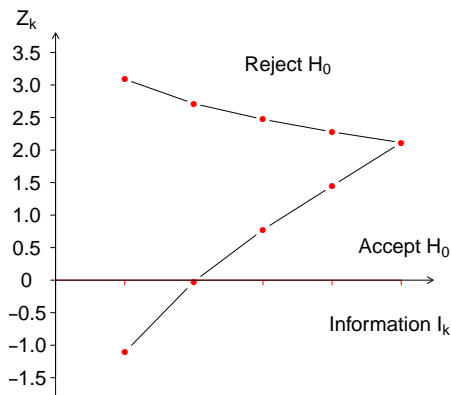
The resulting designs are efficient for this value of \mathcal{I}_{\max} .

The monitoring committee can treat the (non-binding) futility boundary as a guideline, allowing them to consider safety data or secondary endpoints in deciding whether to stop for futility.

An efficacy boundary and non-binding futility boundary

A test with 5 planned analyses, type I error probability $\alpha = 0.025$, power 0.9 if $\theta = \delta = 1$, and type I and II error spending functions

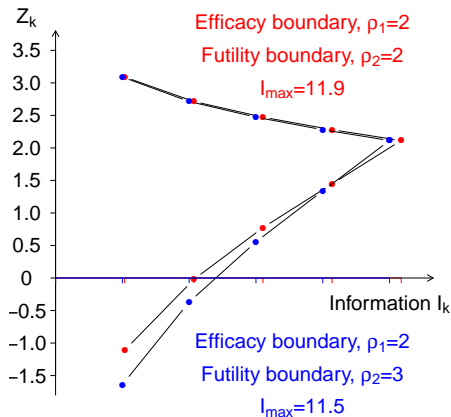
$$f(\mathcal{I}) = \min\{(\mathcal{I}/\mathcal{I}_{\max})^2, 1\} \alpha, \quad g(\mathcal{I}) = \min\{(\mathcal{I}/\mathcal{I}_{\max})^2, 1\} \beta.$$



An efficacy boundary and non-binding futility boundary

Contrast: A test with type I and type II error spending functions

$$f(\mathcal{I}) = \min\{(\mathcal{I}/\mathcal{I}_{\max})^2, 1\} \alpha, \quad g(\mathcal{I}) = \min\{(\mathcal{I}/\mathcal{I}_{\max})^3, 1\} \beta.$$



Resources saved by early stopping

Tests with $\alpha = 0.025$, $1 - \beta = 0.9$ and 5 equally spaced analyses.
Values of \mathcal{I}_{\max} and $E_{\theta}(\mathcal{I})$, expressed as a percentage of \mathcal{I}_{fix} .

<i>Design</i>	\mathcal{I}_{\max}	$E_{\theta}(\mathcal{I})$			
		$\theta=0$	$\theta=0.5$	$\theta=1.0$	$\theta=1.5$
E: $\rho_1 = 2$ only	106	105.2	96.7	70.5	46.8
E: $\rho_1 = 2$, F: $\rho_2 = 2$	113	59.2	80.9	70.6	48.2
E: $\rho_1 = 2$, F: $\rho_2 = 3$	109	64.1	83.0	70.1	47.5

E: Efficacy boundary, F: Non-binding futility boundary

With $\rho_2 = 3$, the maximum sample size is smaller and expected sample size at low values of θ is larger.

The number of analyses and design parameters ρ_1 and ρ_2 can be chosen to give acceptable values of \mathcal{I}_{\max} and $E_{\theta}(\mathcal{I})$.

Conclusions

1. It is common practice for trials to take an informal approach to stopping for futility.

Decision making is often guided by conditional power calculations.

Just how one should use “conditional power” or “predictive power” in deciding whether to stop a trial is not very clear.

2. We can create a group sequential design with a non-binding futility boundary.

The ρ -family of error spending designs provides a simple way to define efficient procedures.

The monitoring committee should be guided by this futility boundary — but they can still use their discretion in deciding whether to stop the trial.

Barber, S. and Jennison, C. (2002). Optimal asymmetric one-sided group sequential tests. *Biometrika*, **89**, 49–60.

Jennison, C., and Turnbull, B. W. (2013). Interim monitoring of clinical trials: decision theory, dynamic programming and optimal stopping. *Kuwait Journal of Science*, **40**, 43–59.

Öhrn, C.F. (2011). *Group Sequential and Adaptive Methods — Topics with Applications to Clinical Trials*. Chapter 4: Group sequential designs with non-binding futility boundaries. PhD thesis, University of Bath.

Rosner, G.L. and Tsiatis, A.A. (1988). The impact that group sequential tests would have made on ECOG clinical trials. *Statistics in Medicine*, **8**, 505–516.

Resources saved by early stopping: Designs with 2 analyses

Tests with $\alpha = 0.025$, $1 - \beta = 0.9$ and **2** equally spaced analyses.

Values of \mathcal{I}_{\max} and $E_{\theta}(\mathcal{I})$, expressed as a percentage of \mathcal{I}_{fix} .

<i>Design</i>	\mathcal{I}_{\max}	$E_{\theta}(\mathcal{I})$			
		$\theta=0$	$\theta=0.5$	$\theta=1.0$	$\theta=1.5$
E: $\rho_1 = 2$ only	103	102.2	97.9	80.5	59.6
E: $\rho_1 = 2$, F: $\rho_2 = 2$	106	70.7	89.2	80.8	60.7
E: $\rho_1 = 2$, F: $\rho_2 = 3$	104	75.5	91.5	80.4	60.0

E: Efficacy boundary, F: Non-binding futility boundary

Resources saved by early stopping: Designs with 3 analyses

Tests with $\alpha = 0.025$, $1 - \beta = 0.9$ and **3** equally spaced analyses.

Values of \mathcal{I}_{\max} and $E_{\theta}(\mathcal{I})$, expressed as a percentage of \mathcal{I}_{fix} .

<i>Design</i>	\mathcal{I}_{\max}	$E_{\theta}(\mathcal{I})$			
		$\theta=0$	$\theta=0.5$	$\theta=1.0$	$\theta=1.5$
E: $\rho_1 = 2$ only	104	103.6	97.1	75.0	52.3
E: $\rho_1 = 2$, F: $\rho_2 = 2$	109	64.4	84.8	75.3	53.5
E: $\rho_1 = 2$, F: $\rho_2 = 3$	106	69.5	86.9	74.8	52.8

E: Efficacy boundary, F: Non-binding futility boundary