

Towards Structured Use of Bayesian Sequential Monitoring in Clinical Trials

Evan Kwiatkowski[†], Eugenio Andraca-Carrera[‡],
Mat Soukup[‡], Matthew A. Psioda^{†*}

[†] Department of Biostatistics, University of North Carolina,
McGavran-Greenberg Hall, CB#7420,
Chapel Hill, North Carolina, U.S.A.

[‡] Division of Biometrics VII, Office of Biostatistics
Center for Drug Evaluation and Research,
US Food and Drug Administration,
Silver Spring, Maryland, USA

September 6, 2019

Abstract

The text of your abstract. 200 or fewer words.

Keywords: 3 to 6 keywords, that do not appear in the title

*The authors gratefully acknowledge *please remember to list all relevant funding sources in the unblinded version*

1 Introduction

Things to discuss:

- 21st Century Cures Act (MATT)
- PDUFA VI reauthorization (MATT)
- Expansive work already done on sequential monitoring (EVAN – draft on 6/21)
- Our majors contribution (EVAN – as early as possible in introduction without having the flow appear weird – draft on 6/21)
- Outline for the remaining section of the paper (EVAN – draft on 6/21)

The theoretical foundations for the Bayesian clinical trials has been long established Cornfield (1966*a*) Cornfield (1966*b*) Neyman & Greenhouse (1967). These methods were not widely used in practice until a comprehensive framework for interpretation of results was developed through specifying prior distributions that were naturally and intuitively related to the research objectives (e.g. skeptical and enthusiastic priors) Freedman & Spiegelhalter (1989) Freedman & Spiegelhalter (1992) Spiegelhalter et al. (1993) Spiegelhalter et al. (1994) Fayers et al. (1997). (*Rewrite paragraph.*)

There is still potential for further utilization of Bayesian methods in the clinical trial setting. While the framework for interpretation of Bayesian clincial trials is well developed, the details of specifying prior distributions in a natural and intuitive way is lacking. This paper presents a structured or default way to determine prior distributions based on the trial design. Our major contribution is to present methods for the default or automatic selection of prior distributions in a way that is applicable to a wide array of clinical trial designs.

1. Bayesian methodology is widely developed.
2. It has been applied (cite).
3. The current perspective is that Bayesian methodology is only valid when Frequentist methods are insufficient, including where enrollment is challenging (rare diseases, pediatric studies)
4. Our contribution is to show that Bayesian methods are applicable to all clinical trials. This is shown by highlighting their improved interpretation and showing their use in varied and complicated situations.

2 Methods

As you introduce ideas that come from or extend other ideas in the literature, cite the relevant literature.

2.1 Monitoring versus Estimation Priors

2.1.1 Bayesian hypothesis testing based on posterior probabilities

The Bayesian paradigm provides direct inference on a parameter of interest through specification of a model for the data generating mechanism and prior distributions for unknown quantities. Let \mathbf{D} be a random variable representing the data collected in the trial with density $p(\mathbf{D}|\theta, \psi)$ where θ and ψ are the unknown quantities. Let θ be the parameter of interest and ψ be the unknown quantities that are not of primary importance (nuisance parameters). Define the sample spaces for the unknown quantities as $\theta \in \Theta$ and $\psi \in \Psi$.

Suppose the hypothesis for the trial is $H_0 : \theta \in \Theta_0$ versus $H_1 : \theta \in \Theta_1$. These hypotheses are judged based on posterior probabilities of θ by evaluating its marginal likelihood

$$P(\theta \in \Theta_i | \mathbf{D}) = \int_{\Theta_i} p(\theta | \mathbf{D}) d\theta \text{ for } i \in \{0, 1\},$$

where $p(\theta | \mathbf{D}) = \int_{\Psi} p(\theta, \psi | \mathbf{D}) d\psi$ is marginalized over the nuisance parameters.

2.1.2 Prior elicitation

It has been said that “the purpose of a trial is to collect data that bring to conclusive consensus at termination opinions that had been diverse and indecisive at the *outset*” (Kass and Greenhouse (1989), emphasis added). These opinions manifest as priors $\pi(\theta, \psi)$ for which their relation to $P(\theta \in \Theta_i | \pi(\theta, \psi))$ $i \in \{0, 1\}$ is examined. Note this quantity does not depend on the data \mathbf{D} and therefore reflects a-priori opinion.

The posterior distribution of θ depends on the choice of prior distribution $\pi(\theta, \psi)$ since $p(\theta, \psi | \mathbf{D}) = p(\mathbf{D} | \theta, \psi) \pi(\theta, \psi) / p(\mathbf{D})$ by Bayes rule. The specification of the prior distribution depends on the research objective. An *inference prior* is a prior that is used when the research objective is to make final analysis after data collection is complete. A *monitoring prior* is a prior that is used when the research objective is to see if there is a persuasive result based in the interim data. Stopping for efficacy is ceasing enrollment due to a promising interim result (one that is consistent with H_1 , and stopping for futility is ceasing enrollment due to a discouraging interim result (one that is consistent with H_0).

Define $\delta \in (0, 1)$ as a threshold for *a compelling level of evidence* as it relates to θ . We say that an individual is “all but convinced” that H_i is true given the observed data if $P(\theta \in \Theta_i | \mathbf{D}) \geq \delta$ for $i \in \{0, 1\}$. The quantity $1 - \delta$ reflects *residual uncertainty* of H_i being true relative to the competing hypothesis.

A enthusiastic prior is an informative prior that gives preference to H_1 such that it is “all but convinced” that H_1 is true a-priori. This prior $\pi_E(\theta, \psi) \equiv \pi_E$ has the property that $P(\theta \in \Theta_1 | \pi_E) \geq \delta$ (equivalently $P(\theta \in \Theta_0 | \pi_E) < 1 - \delta$). The choice of $\delta \in (0, 1)$ is motivated by *a compelling level of evidence* as it relates to θ , although in this setting the “evidence” reflects a theoretical opinion rather than empirical judgement. For example, if $\delta = 0.95$, then this choice of enthusiastic prior places 95% prior probability that $\theta \in \Theta_1$.

A skeptical prior is an informative prior that does not give strong preference to H_1 . This prior $\pi_S(\theta, \psi) \equiv \pi_S$ could have the property that $P(\theta \in \Theta_0 | \pi_S) \geq \delta$, in which case it is “all but convinced” that H_0 is true a-prior, however, this demonstrates such an extreme disbelief in the possibility of a positive effect that conducting the trial at all would be viewed as dubious. Consider a region $\Theta_A \subset \Theta_1$ that demonstrates a sizeable positive effect. The skeptical prior is then constructed such that $P(\theta \in \Theta_A | \pi_S) < \delta$.

2.1.3 Sequential monitoring

The use of monitoring based on changing the opinion of skeptical and enthusiastic priors has been described as overcoming a handicap (Freedman & Spiegelhalter (1989)) and providing a brake (Fayers et al. (1997)) on the premature termination of trials, or constructing “an adversary who will need to be disillusioned by the data to stop further experimentation” (Spiegelhalter et al. (1994)). Early termination of enrollment is appropriate if diverse prior opinions about θ would be in agreement given the interim data (e.g. the skeptical and enthusiastic person reach the same conclusion).

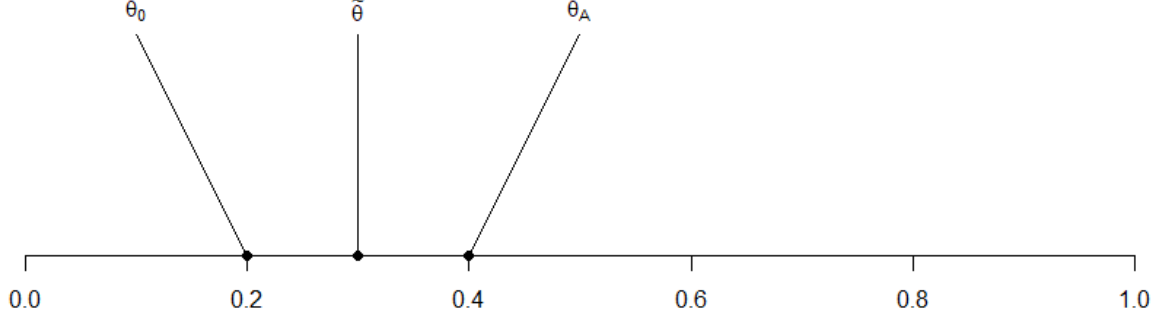


Figure 1:

Promising interim result

In order for interim evidence showing H_1 is true to be persuasive, it has to cause the skeptic, who initially held that $P(\theta \in \tilde{\Theta} \subset \Theta_1 | \pi_S) < \delta$, to believe that $P(\theta \in \Theta_1 | \pi_S) \geq \delta$.

Disillusioning interim result

In order for interim evidence showing H_1 is false to be persuasive, it has to cause the enthusiast, who initially held that $P(\theta \in \Theta_1 | \pi_E) \geq \delta$, to believe that $P(\theta \in \tilde{\Theta} \subset \Theta_1 | \pi_E) < \delta$.

Example

Consider the hypothesis $H_0 : \theta \leq \theta_0$ vs. $H_1 : \theta > \theta_0$. The skeptic initially held that $P(\theta > \theta_A | \pi_S) < \delta$ and a promising interim result would be $P(\theta > \theta_0 | \pi_S) \geq \delta$. The enthusiast initially held that $P(\theta > \theta_0 | \pi_E) \geq \delta$ and a disillusioning interim result would be $P(\theta > \tilde{\theta} | \pi_E) < \delta$.

Probability of Success

As an alternative strategy to futility analysis, one can monitor the probability of success (POS) for the trial. The probability of getting a convincing result at the end of the trail can be computed using the interim data. Let $p(\theta|\mathbf{D}, \pi_I)$ denote the posterior distribution for θ based on the inference prior π_I and the current data \mathbf{D} . Let ξ denote the POS which is given as follows:

$$\begin{aligned}\xi &= P[\mathbf{D}_1 \in \mathbb{R}^{dim(\mathbf{D}_1)} | P(\theta \in \Theta_1 | \mathbf{D}_1, \mathbf{D}, \pi_I) \geq \delta] \\ &= E[1\{P(\theta \in \Theta_1 | \mathbf{D}_1, \mathbf{D}, \pi_I) \geq \delta\}]\end{aligned}$$

where the expectation is taken with respect to the posterior predictive distribution $p(\mathbf{D}_1)$ for future data \mathbf{D}_1 (which includes subjects yet to enroll):

$$p(\mathbf{D}_1) = \int p(\mathbf{D}_1|\theta) \cdot \pi(\theta|\mathbf{D})d\theta.$$

One may stop the enrollment if ξ is sufficiently small (i.e. $\xi < 0.05$).

2.1.4 Final inference

Final inference on the parameter of interest is made once all data has been collected. Enrollment was either stopped based on a persuasive interim result or based on the maximum sample size. An inference prior $\pi_I(\theta, \psi) \equiv \pi_I$ is often less divisive than the skeptical and enthaustic priors, and can be viewed as a balance of the more divisive opinions. We propose use of a mixture prior constructed from the monitoring process as the inference prior:

$$\pi_I = \omega \cdot \pi_S + (1 - \omega) \cdot \pi_E$$

for $\omega \in [0, 1]$. Choosing $\omega = 1/2$ for an equal mixture of π_S and π_E corresponds to an inference prior that equally weights the skeptical and enthuastic opinions. Define

$p(\mathbf{D}|\pi(\theta, \psi)) = \int p(\mathbf{D}|\theta)\pi(\theta, \psi)d(\theta, \psi)$ to be the marginal likelihood for the data given the prior $\pi(\theta, \psi)$. Choosing ω based on posterior model probabilities of the null and alternative hypotheses yields $\omega = p(\mathbf{D}|\pi_S)/(p(\mathbf{D}|\pi_S) + p(\mathbf{D}|\pi_E))$.

All relevant information about θ can be derived from its marginal posterior distribution with an inference prior (e.g. posterior mean, credible intervals). For example, the posterior mean using the inference prior will be a two-part mixture of the posterior means using the skeptical and enthusiastic priors:

$$E(\theta|\mathbf{D}, \pi_I) = \omega \cdot E(\theta|\mathbf{D}, \pi_S) + (1 - \omega) \cdot E(\theta|\mathbf{D}, \pi_E)$$

3 Examples

3.1 Single-Arm Proof-of-Activity Trial with Binary Endpoint

3.1.1 Model formulation & prior elicitation

Consider a single-arm oncology proof-of-activity trial with a binary endpoint. The data \mathbf{D} are Binomially distributed and the response rate θ is the parameter of interest, with higher values of θ being indicative of proof-of-activity. The formulation is discussed in (cite example) with $\theta_0 = 0.2$, $\tilde{\theta} = 0.3$, $\theta_A = 0.4$.

It is intuitive to center the skeptical and enthusiastic priors around the quantities θ_0 and θ_A respectively, so that $E(\pi_S) = \theta_0$ and $E(\pi_E) = \theta_A$.

Beta priors for θ will be used to provide closed-form expressions of the posterior distributions via Beta-Binomial conjugacy (the posterior distribution $p(\theta|\mathbf{D})$ will be Beta distributed). In particular, let y_1 be the number of successes and y_0 be the number of failures. If the skeptical prior is $\pi_S(\theta) \sim \mathcal{B}(\alpha_S, \beta_S)$ then the associated posterior is

$p(\theta|\mathbf{D}, \pi_S) \sim \mathcal{B}(\alpha_S + y_1, \beta_S + y_0)$. Similarly, if the enthusiastic prior is $\pi_E(\theta) \sim \mathcal{B}(\alpha_E, \beta_E)$ then the associated posterior is $p(\theta|\mathbf{D}, \pi_E) \sim \mathcal{B}(\alpha_E + y_1, \beta_E + y_0)$.

The skeptical prior is Beta distributed with expected value $\theta_0 = 0.2$ and has 4.5% prior probability that $\theta > \theta_A = 0.4$. The enthusiastic prior has expected value $\theta_A = 0.4$ and has 5% prior probability that $\theta < \theta_0 = 0.2$. The inference prior will be at an equal mixture of the skeptical and enthusiastic prior ($\omega = 0.5$).

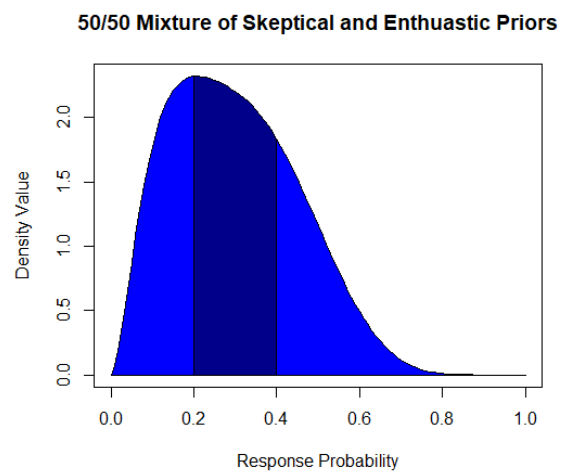
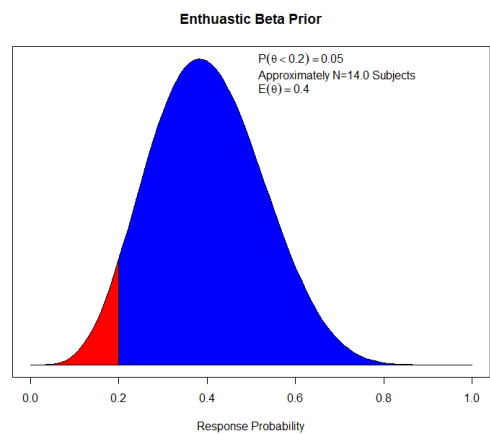
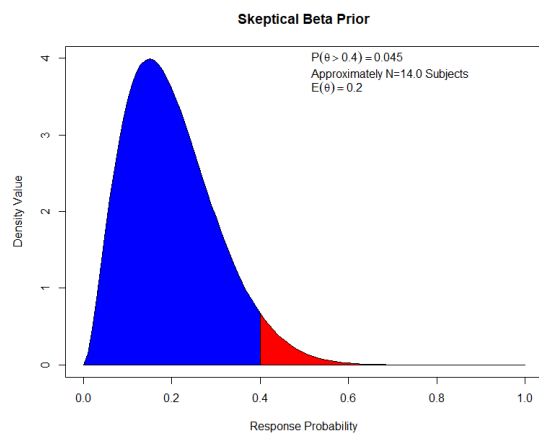


Figure 2: (a) Skeptical prior (b) Enthusiastic prior (c) 50/50 mixture of skeptical and enthusiastic prior

3.1.2 Sequential monitoring

The trial will proceed until one of the following three conditions are satisfied:

Efficacy criteria: $P(\theta > 0.20|\mathbf{D}, \pi_S) \geq 0.95$

Futility criteria: $P(\theta \leq 0.30|\mathbf{D}, \pi_E) \geq 0.85$

Maximum sample size: $N = 76$ patient outcomes obtained

The maximum sample size is based on a frequentist design that would have 90% power at 0.35.

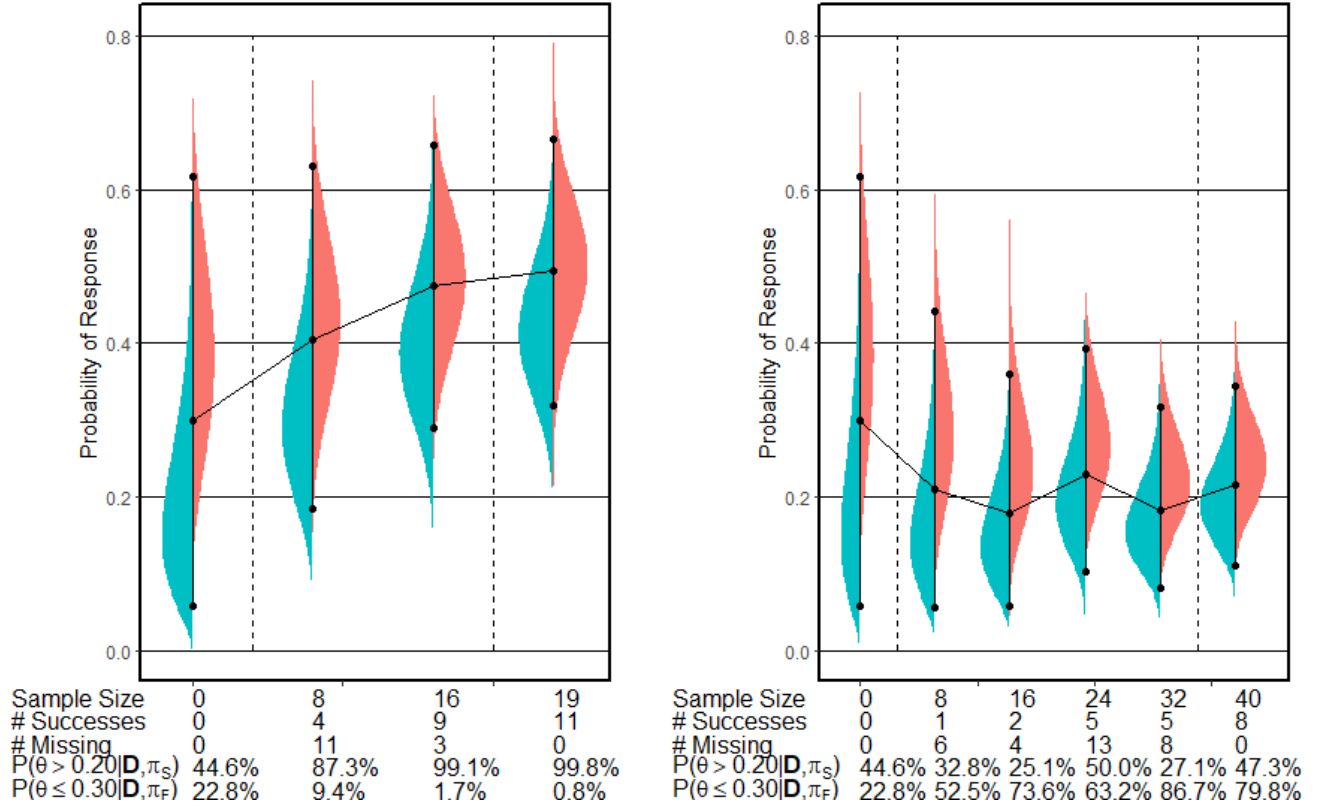


Figure 3: (a) Early stopping for efficacy (b) Early stopping for futility

3.1.3 Example paths

As seen in Figure 3(a), at the second interim analysis the efficacy condition $P(\theta > 0.20 | \mathbf{D}, \pi_S) \geq 0.95$ is satisfied and enrollment is terminated. As shown in Figure 3(b), at the fourth interim analysis the futility condition $P(\theta \leq 0.30 | \mathbf{D}, \pi_F) \geq 0.85$ is satisfied and enrollment is terminated.

3.1.4 Design properties: Results

Assume that the outcomes are ascertained after approximately 4 months of follow-up and 2 patients per month on average are enrolled. An interim analysis will be completed after every 2 subjects complete follow-up. Let EFF be the probability of the trial stopping early for efficacy, FUT be the probability of the trial stopping early for futility, and INC be the probability of reaching the maximum sample size without a conclusive monitoring result. Let SS be the average sample size at the definitive interim analysis (I) and at the end of follow-up (F), let CP be the coverage probability using the mixture prior, and let PM be the posterior mean an inference prior which is a 50/50 mixture of the skeptical and enthusiastic priors.

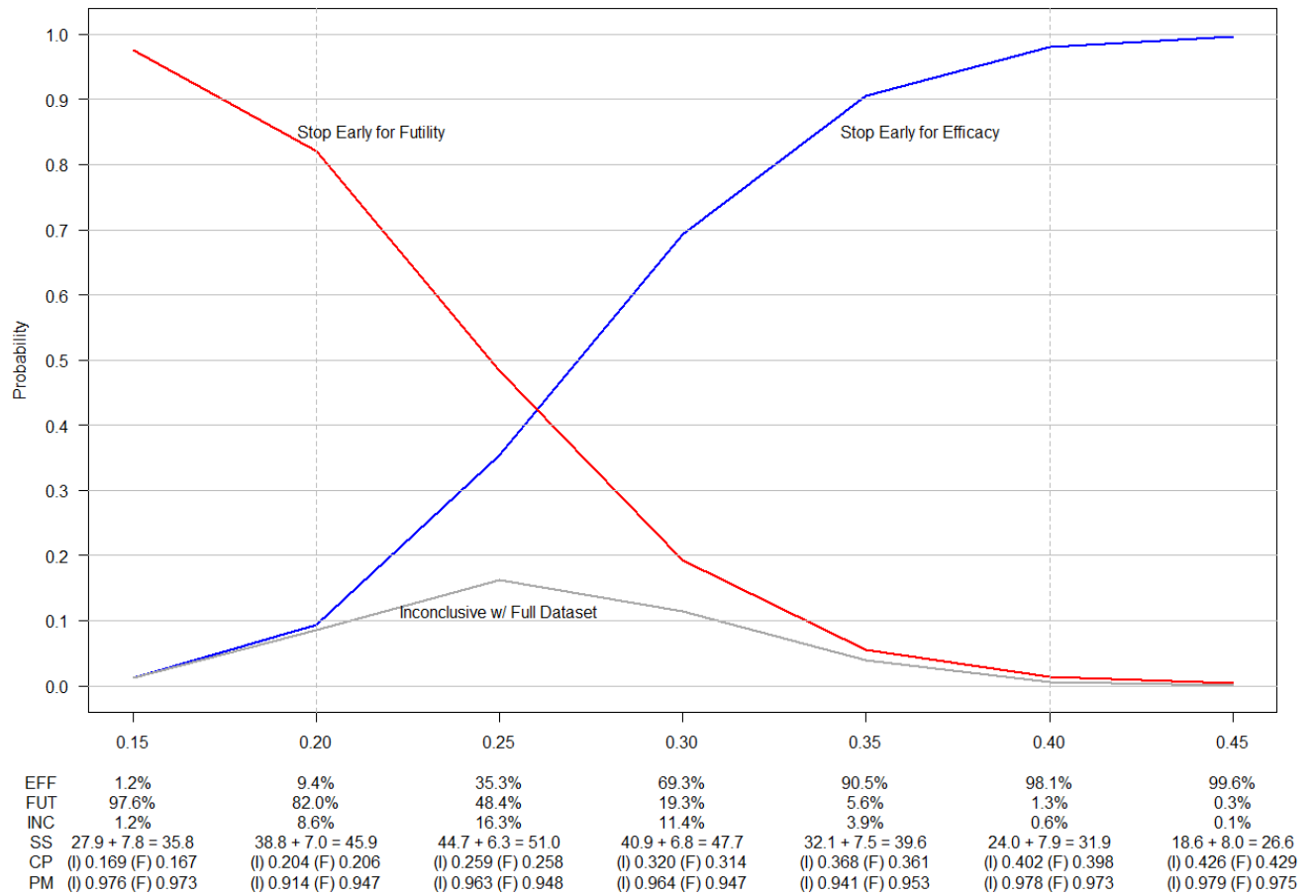


Figure 4:

3.1.5 Agreement between interim and final result

Distribution of final posterior probability given interim stoppage (interim $P(\theta > 0.20|\mathbf{D}, \pi_S) \geq 0.95$) and evidence decrease.

	0.15	0.20	0.25	0.30	0.35	0.40	0.45
Final $P(\theta > 0.20 \mathbf{D}, \pi_S) \geq 0.95$	29.3%	51.0%	64.5%	75.3%	83.2%	89.4%	93.2%
Final $P(\theta > 0.20 \mathbf{D}, \pi_S) < 0.95$	70.7%	49.0%	35.5%	24.7%	16.8%	10.6%	6.8%
Conditional Median	0.91	0.92	0.92	0.93	0.93	0.93	0.93
Conditional 25th percentile	0.87	0.89	0.90	0.91	0.91	0.91	0.91
Conditional 10th percentile	0.83	0.86	0.87	0.88	0.88	0.88	0.88
Conditional 1st percentile	0.67	0.78	0.81	0.81	0.82	0.82	0.82

For example, at a true response rate of $\theta = 0.40$, there is an 89.4% that the threshold for a significant result is maintained after the additional subjects complete follow-up, and in the 10.6% of cases that the evidence decreases, the median posterior probability is 0.93 and only in 10% of cases is the posterior probability lower than 0.88. Thus there is a slight attenuation with respect to the dichotomous threshold, but little change in the posterior probability overall.

3.1.6 Type 1 error rate by the frequency of data monitoring

As expected, the probability of stopping enrollment due to a promising interim trial result and the Type 1 error rate at the final analysis increase with the frequency of interim monitoring, however, the increase is very slight at the final analysis. Regardless of frequency of monitoring there are good Type 1 error rates. Even at the extreme case where an interim analysis is conducted after every outcome, the probability of stopping at the interim due to a promising result when the true response is at the null level is only 0.108 (about double the

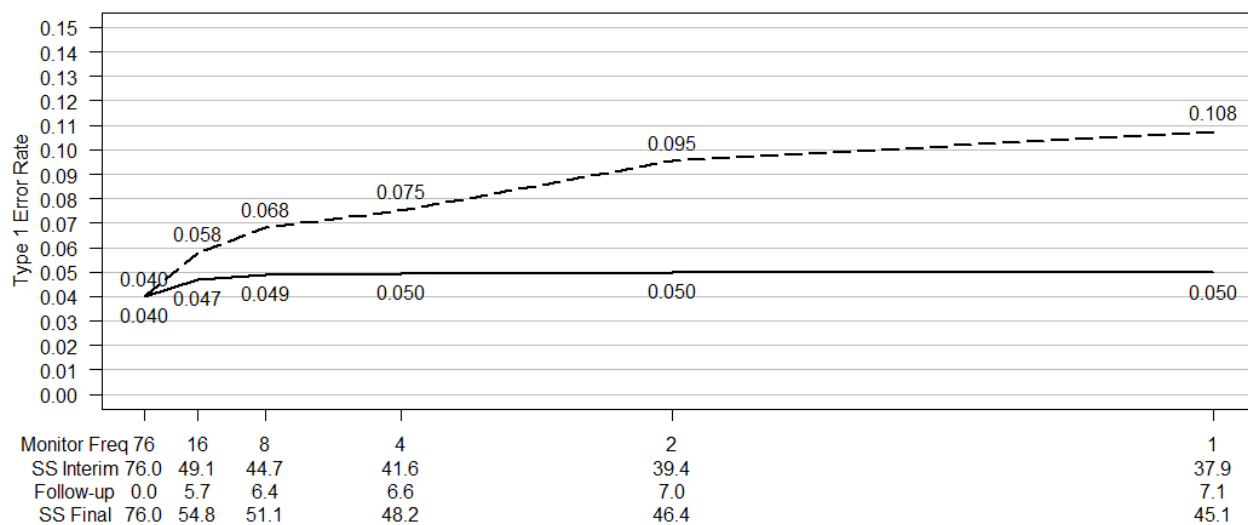


Figure 5: Type 1 error rate depending on frequency of sequential monitoring

nominal rate), and even in this situation the Type 1 error rate once follow-up is complete does not exceed 0.05. Thus Bayesian sequential monitoring has good frequentist properties even with frequent interim analyses.

3.1.7 Type 1 error rate depending on enrollment schemes

Consider the same trial but with a longer follow-up length of 8 months rather than 4 months.

Comparison of Efficacy Stopping, Type 1 Error Rate, and Sample Size by follow-up length and frequency of sequential monitoring

Mon Freq	Eff. Stopping		T1E Final		SS Final		% Ongoing	
	4 mon	8 mon	4 mon	8 mon	4 mon	8 mon	4 mon	8 mon
76	0.040	0.039	0.040	0.039	76.0	76.0	0.0%	0.0%
16	0.058	0.056	0.047	0.042	54.8	60.0	10.4%	18.1%
8	0.068	0.067	0.049	0.043	51.1	56.7	12.5%	21.3%
4	0.075	0.075	0.050	0.043	48.2	54.1	13.7%	23.4%
2	0.095	0.094	0.050	0.043	46.4	52.8	15.1%	25.5%
1	0.108	0.107	0.050	0.043	45.1	51.7	15.7%	26.8%

Note that the probability of efficacy stopping and Type 1 error rate increase monotonically for both specifications of follow-up length. The Type 1 error rate is lower for the 8-month follow-up design since there are more subjects in the final sample size.

4 Robustness of parameterizations of monitoring priors

The spike-slab and flattened priors are both 2-part mixtures of Beta distributions.

- Test the hypothesis that power curves will be identical and what will change is the expected sample sizes and that varying the decay rate will affect bias not Type I/II error
- Actually, having one prior “default” and the other spike-slab changes the monitoring result substantially. When using the spike-slab skeptical prior and the default enthusiastic prior, the determination of efficacy at the interim is made less frequently than when using the default skeptical prior. Similarly, when using the spike-slab enthusiastic prior and the default skeptical prior, the determination of futility is made less frequently than when using the default enthusiastic prior.

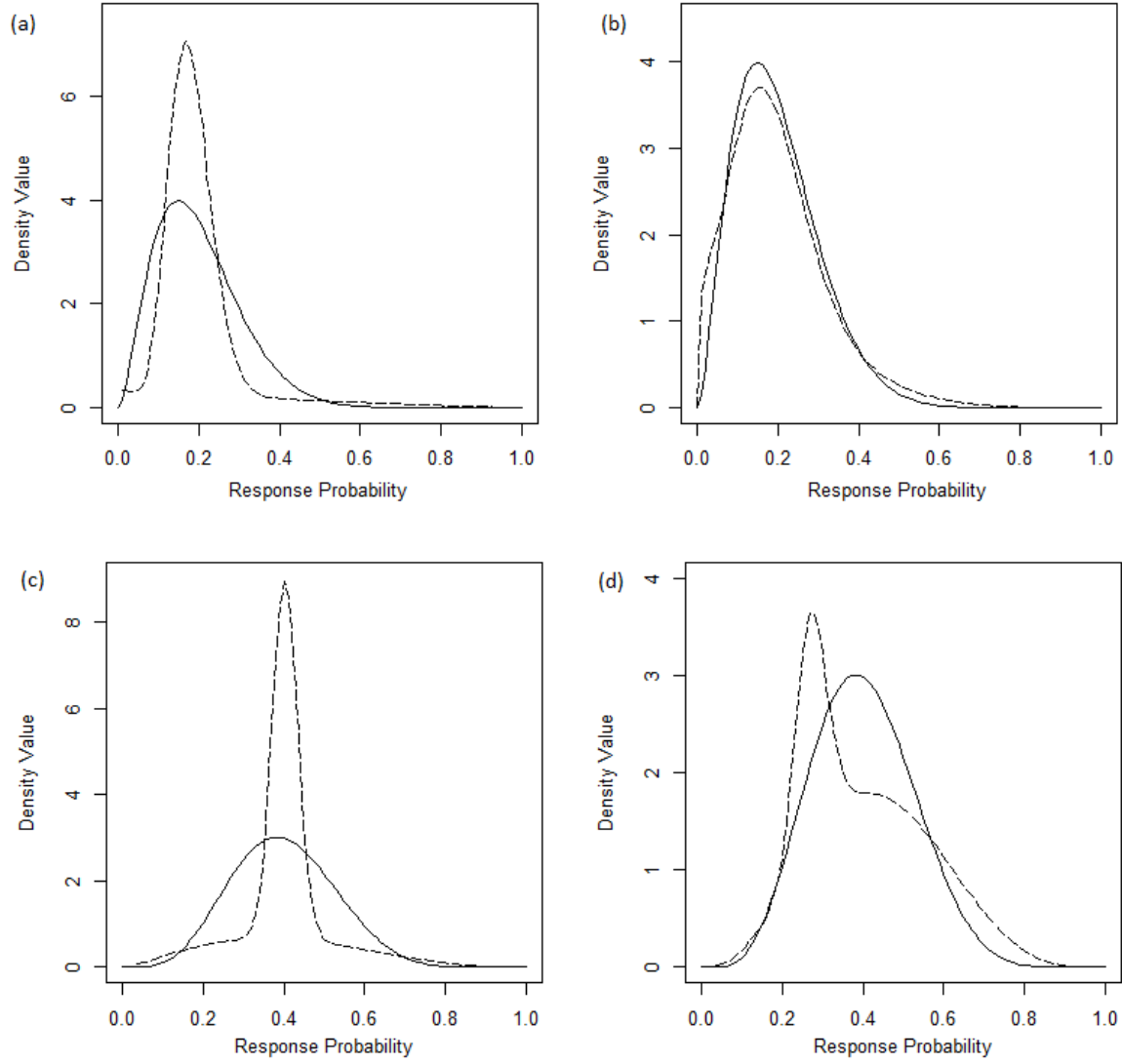


Figure 6: (a) Spike-slab skeptical prior (b) Flattened skeptical prior (c) Spike-slab enthusiastic prior (d) Flattened enthusiastic prior

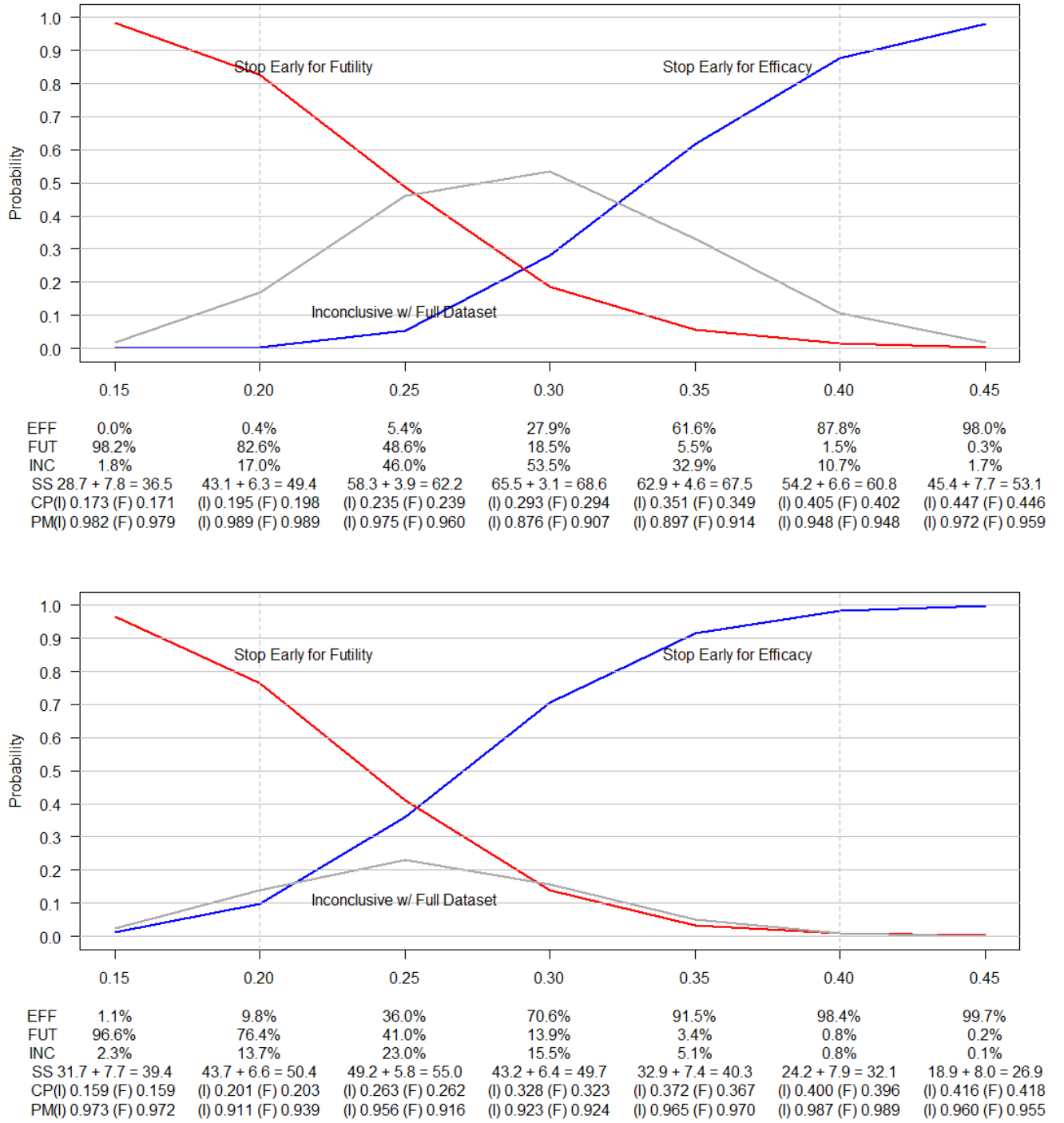


Figure 7: (a) Spike-slab skeptical prior and default enthusiastic prior (b) Default skeptical prior and spike-slab enthusiastic prior

4.1 Parallel Two-Group Superiority Trial /w Continuous Binary Endpoint

Interesting because prior is on risk difference $[-1,1]$ while also being non-informative on control group. Will need numerical integration to evaluate posteriors.

4.2 Three-Arm, Placebo Controlled Non-Inferiority Trial w/ Continuous Endpoint

$$P \rightarrow \beta_0 \text{ (placebo)}$$

$$C \rightarrow \beta_0 + \beta_1 \text{ (control)}$$

$$A \rightarrow \beta_0 + \beta_1 + \beta_2 \text{ (active)}$$

$$H_0 : \beta_2 - \delta\beta_1 \leq 0$$

Parameters of interest (β_1, β_2) , nuisance parameters (β_0, σ^2) .

Need priors $\pi(\beta_0), \pi(\beta_1), \pi(\beta_2|\beta_1)$.

Will use MCMC to evaluate posteriors.

5 Discussion – (MATT/EVAN)

Q: Why not reverse engineer priors to have exact Type 1 error properties?

A: This would basically be a frequentist method, in that the design would have to be adhered to exactly (including number and timing of data monitoring). Philosophically, designing a Bayesian study that requires rigid monitoring rules loses the advantages of Bayes from the likelihood principle.

SUPPLEMENTARY MATERIAL

6 Beta Priors

Beta priors for θ will be used to provide closed-form expressions of the posterior distributions via Beta-Binomial conjugacy (the posterior distribution $p(\theta|\mathbf{D})$ will be Beta distributed). The Beta distribution has two shape parameters. These parameters can be determined uniquely by specifying the desired mean and variance of the distribution. The variance for the skeptical and enthusiastic priors is then uniquely determined through by the choice of threshold δ . In particular, let $\pi_S(\theta) \sim \mathcal{B}(\alpha, \beta)$ be Beta distributed with shape parameters (α, β) . There is a single choice of (α, β) such that:

$$\theta_0 = E(\pi_S) = \int_{\Theta} \pi_S(\theta) d\theta = \frac{\alpha}{\alpha + \beta} \text{ and } \delta = \int_{\Theta_0} \pi_S(\theta) d\theta = \int_0^{\theta_0} \frac{\theta^{\alpha-1}(1-\theta)^{\beta-1}}{B(\alpha, \beta)} d\theta$$

where $B(\alpha, \beta)$ is the Beta function.

Alternatively, the variance could be determined by specifying a desired quantile of the prior distribution which would then be reflected in δ . Then there is a single choice of (α, β) such that

$$\theta_0 = E(\pi_S) = \int_{\Theta} \pi_S(\theta) d\theta = \frac{\alpha}{\alpha + \beta} \text{ and } \lambda = \int_{\theta_A}^1 \pi_S(\theta) d\theta = \int_{\theta_A}^1 \frac{\theta^{\alpha-1}(1-\theta)^{\beta-1}}{B(\alpha, \beta)} d\theta,$$

in which case $\delta = \int_{\Theta_0} \pi_S(\theta) d\theta$ is a deterministic quantity.

7 BibTeX

References

Cornfield, J. (1966*a*), ‘A Bayesian Test of Some Classical Hypotheses, with Applications to Sequential Clinical Trials’, *Journal of the American Statistical Association* **61**(315), 577.
URL: <https://www.jstor.org/stable/2282772?origin=crossref>

Cornfield, J. (1966*b*), ‘Sequential Trials, Sequential Analysis and the Likelihood Principle’, *The American Statistician* **20**(2), 18.
URL: <https://www.jstor.org/stable/2682711?origin=crossref>

Fayers, P. M., Ashby, D. & Parmar, M. K. B. (1997), ‘Tutorial in Biostatistics: Bayesian Data Monitoring in Clinical Trials’, *Statistics in Medicine* **16**(12), 1413–1430.
URL: <http://doi.wiley.com/10.1002/%28SICI%291097-0258%2819970630%2916%3A12%3C1413%3A%3E1.0.CO%3B2-U>

Freedman, L. S. & Spiegelhalter, D. J. (1989), ‘Comparison of Bayesian with group sequential methods for monitoring clinical trials’, *Controlled Clinical Trials* **10**(4), 357–367.
URL: <https://www.sciencedirect.com/science/article/pii/0197245689900019?via%3Dihub>

Freedman, L. S. & Spiegelhalter, D. J. (1992), ‘Application of bayesian statistics to decision making during a clinical trial’, *Statistics in Medicine* **11**(1), 23–35.
URL: <http://doi.wiley.com/10.1002/sim.4780110105>

Neyman, J. & Greenhouse, S. W. (1967), *Proceedings of the Berkeley Symposium on Mathematical Statistics and Probability.*, University of California Press.
URL: <https://projecteuclid.org/euclid.bsmsp/1200513830>

Spiegelhalter, D. J., Freedman, L. S. & Parmar, M. K. B. (1993), ‘Applying Bayesian ideas in drug development and clinical trials’, *Statistics in Medicine* **12**(15-16), 1501–1511.

URL: <http://doi.wiley.com/10.1002/sim.4780121516>

Spiegelhalter, D. J., Freedman, L. S. & Parmar, M. K. B. (1994), ‘Bayesian Approaches to Randomized Trials’, *Journal of the Royal Statistical Society. Series A (Statistics in Society)* **157**(3), 357.

URL: <https://www.jstor.org/stable/10.2307/2983527?origin=crossref>