



A Utility-Based Design for Randomized Comparative Trials with Ordinal Outcomes and Prognostic Subgroups

Thomas A. Murray ^{1,*} Ying Yuan ^{2,**} Peter F. Thall,^{2,***} Joan H. Elizondo,³
and Wayne L. Hofstetter⁴

¹Division of Biostatistics, University of Minnesota, Minneapolis, Minnesota, U.S.A.

²Department of Biostatistics, The University of Texas M.D. Anderson Cancer Center, Houston, Texas, U.S.A.

³Department of Clinical Nutrition, The University of Texas M.D. Anderson Cancer Center, Houston, Texas, U.S.A.

⁴Department of Thoracic and Cardiovascular Surgery, The University of Texas M.D. Anderson Cancer Center, Houston, Texas, U.S.A.

**email:* murra484@umn.edu

***email:* yyuan@mdanderson.org

****email:* rex@mdanderson.org

SUMMARY. A design is proposed for randomized comparative trials with ordinal outcomes and prognostic subgroups. The design accounts for patient heterogeneity by allowing possibly different comparative conclusions within subgroups. The comparative testing criterion is based on utilities for the levels of the ordinal outcome and a Bayesian probability model. Designs based on two alternative models that include treatment-subgroup interactions are considered, the proportional odds model and a non-proportional odds model with a hierarchical prior that shrinks toward the proportional odds model. A third design that assumes homogeneity and ignores possible treatment-subgroup interactions also is considered. The three approaches are applied to construct group sequential designs for a trial of nutritional prehabilitation versus standard of care for esophageal cancer patients undergoing chemoradiation and surgery, including both untreated patients and salvage patients whose disease has recurred following previous therapy. A simulation study is presented that compares the three designs, including evaluation of within-subgroup type I and II error probabilities under a variety of scenarios including different combinations of treatment-subgroup interactions.

KEY WORDS: Group sequential; Hierarchical model; Non-proportional odds; Ordinal response; Precision medicine.

1. Background

This article describes a design for a small single-center randomized controlled clinical trial to evaluate the effectiveness of nutritional prehabilitation (Nuprehab) for esophageal cancer patients who undergo esophageal resection preceded and followed by chemoradiation therapy. Common postoperative morbidities for patients who undergo esophageal resection include anastomotic leak and stricture, chylothorax, delayed emptying, or dumping syndrome, pulmonary complications such as pneumonia, and cardiac complications such as atrial fibrillation (Parekh and Iannettoni, 2007; Chen, 2014). The Nuprehab is given prior to surgery as well as seven days after surgery, with the aim to achieve oral immunomodulation with an L-arginine based enteral formula. The motivation for the trial is the hypothesis that providing patients with Nuprehab may reduce the incidence of postoperative morbidity and mortality via nutritional supplementation (see, e.g., Braga et al., 2002; Waitzberg et al., 2006).

Patients will be randomized to receive either Nuprehab or control, which is the standard of care. All patients will be monitored for morbidity and mortality for 30 days following

their surgery. The trial's primary outcome is Clavien–Dindo postoperative morbidity (POM) score (Clavien et al., 1992; Dindo et al., 2004; Clavien et al., 2009), which is ordinal with six levels: 0 = normal recovery, 1 = minor complication, 2 = complication requiring pharmaceutical intervention, 3 = complication requiring surgical, endoscopic or radiological intervention, 4 = life-threatening complication requiring intensive care, and 5 = death. The worst POM score during 30 days post surgery will be recorded. The trial will enroll approximately 60% primary and 40% salvage patients. Primary patients are treatment naive, whereas salvage patients have been treated previously with chemoradiation therapy, but not surgery, and their disease has recurred. Salvage patients are expected to have fewer preoperative nutritional deficiencies, but more preoperative comorbidities and worse prognosis. Consequently, it is plausible that the efficacy of Nuprehab may differ substantially for primary and salvage patients.

In this article, we describe a design that accounts for the possibility that Nuprehab may be clinically beneficial for one of the subgroups but not the other. This is in sharp contrast with a more traditional “one-size-fits-all” approach that ignores prognostic information and makes one recommendation for all patients about whether Nuprehab is clinically

[Corrections added on September 11, 2018, after first online publication: Acknowledgment section added]

beneficial. We evaluate the design based on its probabilities of recommending Nuprehab to each subgroup in four key scenarios: (i) the “complete null” scenario where Nuprehab does not improve POM scores for patients in either subgroup; (ii) the “partial null” scenario where Nuprehab does not improve POM scores for primary patients, but achieves a targeted benefit for salvage patients; (iii) the “partial null” scenario where Nuprehab achieves a targeted benefit for primary patients, but does not improve POM scores for salvage patients; and (iv) the “complete alternative” scenario where Nuprehab achieves targeted benefits in POM score reduction for both primary and salvage patients. The proposed design addresses the concern that, in the partial null scenarios (ii) and (iii), a one-size-fits-all design will have an unacceptably high (low) probability for recommending Nuprehab to the non-benefiting (benefiting) subgroup.

The proposed design is frequentist in that it is specified to provide specific probabilities of recommending Nuprehab to a non-benefiting subgroup (i.e., type I error) and to a benefiting subgroup for a particular targeted benefit (i.e., power). However, the decision to recommend Nuprehab for a particular subgroup is based on a posterior probability from a Bayesian model. Designs have been proposed that are similar to our design in that they facilitate subgroup specific recommendations, and stopping subsequent enrollment for particular subgroups, see, for example, Brannath et al. (2009), Wang et al. (2009), Rosenblum et al. (2016). These designs do not involve an ordinal outcome, however.

We consider two alternative Bayesian probability models that facilitate subgroup specific recommendations. The first is the proportional odds (PO) cumulative logistic regression model of McCullagh (1980) with a treatment-subgroup interaction parameter. The second is a non-proportional odds (NPO) model with a hierarchical prior that shrinks toward the PO model. Although NPO models have been proposed, see, for example, Peterson and Harrell (1990), Bender and Grouven (1998), Ishwaran (2000), Agresti (2010), as far as we are aware, our formulation of the NPO model is novel, and moreover this is the first proposal to use an NPO model as the basis for comparing treatments in a randomized clinical trial. Guo and Yuan (2017) use the dispersed cumulative probit model of McCullagh (1980), which is a type of NPO model, for personalized dose finding in a phase I/II study of molecularly targeted agents. As we describe below, compared with the PO model, treatment comparison based on the proposed NPO model is more robust but more complex. To obtain a practical design that deals with this complexity, we propose a comparative testing criterion based on elicited numerical utilities of the six POM scores. Our approach may be considered a generalization of the utility-based design proposed by Murray et al. (2016), which does not accommodate prognostic subgroups and uses a Dirichlet–multinomial model.

The remainder of the article is organized as follows. In Section 2, we discuss treatment comparison with ordinal outcomes in general and our utility-based comparative criterion in particular. In Section 3, we describe the PO and NPO models. In Section 4, we discuss practical design considerations, including specifying targeted alternatives, analysis and mon-

itoring plan, controlling the probability of committing a type I error, and sample size. In Section 5, we present the results of a simulation study comparing the proposed design based on either the PO or NPO model, and also a more traditional design based on a PO model without a treatment-subgroup interaction parameter. We conclude with a brief discussion in Section 6.

2. Treatment Comparison

Each design compares the efficacy of Nuprehab relative to standard of care using the six-level ordinal POM score. Comparing treatments based on an ordinal outcome is complicated by the fact that, even when the probability of each outcome level is known, it is not always clear whether one treatment is superior to the other. A simple example is a three-level outcome (Good, Intermediate, Poor) where treatment A gives probabilities (0.30, 0.50, 0.20) and treatment B gives probabilities (0.40, 0.30, 0.30). Since B has larger probabilities of both Good and Poor compared to A , it is not clear whether one treatment is superior to the other. Comparing the two treatments requires additional information, such as a quantification of the relative desirabilities of the three possible events.

Accounting for prognostic subgroups further complicates matters. To see this, denote $Y =$ POM score, $P =$ primary, $S =$ salvage, $N =$ Nuprehab, $C =$ control, and

$$\pi_y(\text{Sgp}, \text{Trt}) = \text{Prob}(Y = y | \text{Sgp}, \text{Trt}) \text{ and}$$

$$\pi_y^+(\text{Sgp}, \text{Trt}) = \text{Prob}(Y \leq y | \text{Sgp}, \text{Trt}),$$

for $y = 0, \dots, 5$, $\text{Sgp} \in \{P, S\}$ and $\text{Trt} \in \{N, C\}$. Indexing Sgp by x and Trt by a or a' , if

$$\pi_y^+(x, a) \geq \pi_y^+(x, a'), \text{ for } y = 0, 1, \dots, 4, \text{ and}$$

$$\pi_y^+(x, a) > \pi_y^+(x, a'), \text{ for some } y = 0, 1, \dots, 4,$$

then clearly treatment a is superior to a' for patients in subgroup x . By contrast, if

$$\pi_y^+(x, a) < \pi_y^+(x, a'), \text{ for some } y = 0, 1, \dots, 4, \text{ and}$$

$$\pi_y^+(x, a) > \pi_y^+(x, a'), \text{ for some } y = 0, 1, \dots, 4,$$

then it is not clear whether a is superior to a' for patients in subgroup x . Stated formally, for a particular patient subgroup, if the POM score distributions corresponding to each treatment arm are *stochastically ordered* with strict inequality $\pi_y^+(x, a) > \pi_y^+(x, a')$ for at least one level y , then it is clear which treatment is superior for that subgroup. By contrast, if the POM score distributions corresponding to each treatment arm are not stochastically ordered, then it is not clear whether one treatment is superior to the other for that subgroup.

To provide a criterion for determining whether one treatment is superior, we elicit numerical utilities $U(Y = y)$ for all levels of Y , and compare treatments using mean utilities,

$$\bar{U}(\pi(\text{Sgp}, \text{Trt})) = \sum_{y=0}^5 U(Y = y) \times \pi_y(\text{Sgp}, \text{Trt}).$$

These depend on the subgroup-treatment specific outcome probabilities $\pi(\text{Sgp, Trt})$ and a utility function $U(Y = y)$ that quantifies the desirability of each outcome level. Following, Houede et al. (2010), Thall and Nguyen (2012), and Murray et al. (2016), we elicited $\{U(Y = y), y = 0, \dots, 5\}$ from the trial's Principal Investigator, WH, so that the numerical utilities reflect his familiarity with postoperative complications following esophageal resection. To do this, we first set $U(Y = 0) = 100$ and $U(Y = 5) = 0$, and then asked WH to specify numerical values for the intermediate levels, $y = 1, \dots, 4$, that reflect their desirability relative to the best and worst levels. The numerical values that WH chose are:

$$\begin{aligned} U(Y = 0) &= 100, & U(Y = 1) &= 80, & U(Y = 2) &= 65, \\ U(Y = 3) &= 25, & U(Y = 4) &= 10, & U(Y = 5) &= 0. \end{aligned}$$

These reflect that POM scores ≤ 2 are substantially more desirable than POM scores ≥ 3 . Because a larger mean utility corresponds to better patient outcomes on average, if $\bar{U}(\pi(x, a)) > \bar{U}(\pi(x, a'))$, then treatment a is superior to a' for patients in subgroup x . Therefore, regardless of whether $\pi(x, a)$ and $\pi(x, a')$ are stochastically ordered, the mean utilities provide an unambiguous criterion for comparing treatments.

One important property of the mean utilities is a consequence of the following theorem.

THEOREM 1. *If $\pi(x, a)$ stochastically dominates $\pi(x, a')$, then $\bar{U}(\pi(x, a)) > \bar{U}(\pi(x, a'))$ for all admissible $U(Y)$ such that $U(Y = 0) > U(Y = 1) > \dots > U(Y = 4) > U(Y = 5)$.*

Theorem 1 follows from first-order stochastic dominance (Quirk and Saposnik, 1962); nonetheless, we provide a proof in the supplementary materials. Consequently, when the POM score distributions are stochastically ordered—and thus, it is clear which treatment is superior without appealing to the mean utilities—the proposed utility-based comparison is not sensitive to the elicited numerical values in that a different set of admissible values will result in the same conclusion. By contrast, when the POM score distributions are not stochastically ordered, eliciting numerical values is necessary to determine whether one treatment is superior. The proposed utility-based comparison necessarily is sensitive to the elicited values in that a different set of admissible values may result in a different conclusion.

Since the POM score probabilities are unknown, we learn about these using a Bayesian model with unknown parameter θ and model-based mean utilities $\bar{U}(\pi(\text{Sgp, Trt}; \theta))$. Given interim or final data \mathbf{D} , our comparative testing criterion is as follows. If

$$\text{Prob}\{\bar{U}(\pi(x, a; \theta)) > \bar{U}(\pi(x, a'; \theta)) \mid \mathbf{D}\} > p_{cut},$$

then we declare treatment a is superior to a' for patients in subgroup x . We specify p_{cut} to control subgroup specific type I error probabilities. We describe how to do this in Section 4.

3. Probability Models

During the process of designing the trial, we considered two Bayesian cumulative logistic regression models that both include treatment-subgroup interaction parameters. The first is a PO model, which is a popular regression model for ordinal response variables, see, for example, McCullagh (1980), Walters et al. (2001), Abreu et al. (2008). The restrictive parametric assumption underlying the PO model often is unrealistic, however. The second is a NPO model that relaxes this assumption at the cost of greater model complexity.

3.1. Proportional Odds Model

Denote $\text{logit}(q) = \log\{q/(1 - q)\}$. The PO model that we considered assumes

$$\begin{aligned} \text{logit}\{\pi_y^+(X, A; \alpha_y, \beta)\} &= \alpha_y + \beta_1 X + \beta_2 A + \beta_3 XA, \\ &\text{for } y = 0, \dots, 4, \end{aligned} \tag{1}$$

where $\alpha_0 \leq \dots \leq \alpha_4$, $X = -0.5$ for primary, $X = 0.5$ for salvage, $A = -0.5$ for control and $A = 0.5$ for Nuprehab. This model is parsimonious in that it accounts for all treatment and subgroup effects using three parameters, $\beta = (\beta_1, \beta_2, \beta_3)$.

Let $n_y(X, A)$ denote the number of patients with a POM score equal to y in subgroup X and treatment arm A . We assume observations are mutually independent, so that the likelihood function for the unknown model parameters (α, β) is

$$L(\alpha, \beta \mid \mathbf{D}) = \prod_{y \in \{0, \dots, 5\}} \prod_{X \in \{-0.5, 0.5\}} \prod_{A \in \{-0.5, 0.5\}} \pi_y(X, A; \alpha, \beta)^{n_y(X, A)}, \tag{2}$$

where

$$\begin{aligned} \pi_0(X, A; \alpha, \beta) &= \pi_0^+(X, A; \alpha_0, \beta), \\ \pi_y(X, A; \alpha, \beta) &= \pi_y^+(X, A; \alpha_y, \beta) - \pi_{y-1}^+(X, A; \alpha_{y-1}, \beta), \\ &\text{for } y = 1, \dots, 4, \\ \pi_5(X, A; \alpha, \beta) &= 1 - \pi_4^+(X, A; \alpha_4, \beta), \end{aligned}$$

and $\pi_y^+(X, A; \alpha_y, \beta)$, $y = 0, \dots, 4$, is defined in (1).

We specify the prior distribution for (α, β) such that $p_0(\alpha, \beta) = p_0(\alpha) \times p_0(\beta)$, where

$$p_0(\alpha) = p_0(\alpha_0) \times \prod_{s=1}^4 p_0(\alpha_s \mid \alpha_{s-1})$$

$$\text{and } p_0(\beta) = p_0(\beta_1) \times p_0(\beta_2) \times p_0(\beta_3).$$

The exact prior distributional forms that we assume are

$$\begin{aligned} \alpha_0 &\sim t_5(\alpha_0^*, 2.5), & \alpha_y \mid \alpha_{y-1} &\sim t_5(\alpha_y^*, 2.5)[\alpha_{y-1}, \infty], \\ &\text{for } y = 1, 2, 3, 4, \\ \beta_1 &\sim t_5(\beta_1^*, 2.5), & \beta_2 &\sim t_5(0, 2.5), & \beta_3 &\sim t_5(0, 2.5), \end{aligned} \tag{3}$$

where we write $p_0(\theta)[L, U]$ to denote that $p_0(\theta)$ has support on the interval $[L, U]$. The above prior restricts $\alpha_0 \leq \dots \leq \alpha_4$

Table 1

Elicited prior mean POM score probabilities in the control arm for primary and salvage patients

Subgroup	POM score					
	0	1	2	3	4	5
Primary	0.50	0.20	0.10	0.10	0.05	0.05
Salvage	0.30	0.25	0.10	0.10	0.10	0.15

so that $0 \leq \pi_1^+(X, A; \alpha_1, \beta) \leq \dots \leq \pi_4^+(X, A; \alpha_4, \beta) \leq 1$ for all $X \in \{-.5, .5\}$ and $A \in \{-.5, .5\}$, which is necessary to ensure that the probability model is admissible. Following the recommendations of Ghosh et al. (2017), we specify t-distributions with a scale of 2.5 and five degrees of freedom. This specification places about 90% of the prior probability mass on the range of values within 5 of the prior mean, while the heavy tails do not preclude more extreme values should the data demand this. Because an effect size of 5 corresponds to a shift from 0.01 to 0.99 on the probability domain between subgroups or treatment arms, the proposed prior specification allows the observed data to dominate posterior inference.

We specify non-zero prior means, $\{\alpha_y^*\}_{y=0}$ and β_1^* , to reflect the prior information that WH provided about the POM score probabilities of each subgroup in the control arm, which we report in Table 1. Using $\pi_y^*(x)$ to denote the prior probabilities that WH provided for the subgroup corresponding to $X = x$, and $\pi_y^{+,*}(x) = \sum_{\ell=1}^y \pi_\ell^*(x)$, we set

$$\alpha_y^* = \left[\log \left(\frac{\pi_y^{+,*}(P)}{1 - \pi_y^{+,*}(P)} \right) + \log \left(\frac{\pi_y^{+,*}(S)}{1 - \pi_y^{+,*}(S)} \right) \right] / 2,$$

for $y = 0, \dots, 4$, and

$$\beta_1^* = \left[\sum_{y=0}^4 \log \left(\frac{\pi_y^{+,*}(P)}{1 - \pi_y^{+,*}(P)} \right) - \log \left(\frac{\pi_y^{+,*}(S)}{1 - \pi_y^{+,*}(S)} \right) \right] / 5.$$

Because we set the prior means for β_2 and β_3 equal to zero, *a priori* $\bar{U}(\pi(P, N)) = \bar{U}(\pi(P, C))$ and $\bar{U}(\pi(S, N)) = \bar{U}(\pi(S, C))$. Therefore, *a posteriori* a non-zero mean utility difference between treatment arms in either subgroup will reflect the observed data, and not the prior.

The PO model assumes that the regression coefficients, and thus the log-odds ratios, do not differ with the level of the response variable. This is a strong parametric assumption that often is unrealistic in practice, including the present context. The PO model likely is popular since it facilitates treatment comparison in that a utility function need not be elicited from the clinician(s). To see this, note that for $y = 0, \dots, 4$ and $X = x$,

$$\pi_y^+(x, N; \alpha_y, \beta) = \frac{\pi_y^+(x, C; \alpha_y, \beta) \exp\{\beta_2 + \beta_3 x\}}{1 - \pi_y^+(x, C; \alpha_y, \beta) + \pi_y^+(x, C; \alpha_y, \beta) \exp\{\beta_2 + \beta_3 x\}},$$

$\pi_y^+(x, N; \alpha_y, \beta)$ is monotonically increasing in $(\beta_1 + \beta_3 x)$ such that $\pi_y^+(x, N; \alpha_y, \beta) = \pi_y^+(x, C; \alpha, \beta)$ when $(\beta_1 + \beta_3 x) =$

0. Therefore, for any $U(Y = 0) < \dots < U(Y = 5)$,

$$\text{if } (\beta_1 + \beta_3 x) > 0, \text{ then } \bar{U}(\pi(x, N; \beta, \alpha)) > \bar{U}(\pi(x, C; \beta, \alpha)),$$

and conversely,

$$\text{if } (\beta_1 + \beta_3 x) < 0, \text{ then } \bar{U}(\pi(x, N; \beta, \alpha)) < \bar{U}(\pi(x, C; \beta, \alpha)).$$

Consequently, when posterior inference is based on the PO model defined in (1), the utility function is superfluous. However, when the actual response distributions do not satisfy the PO assumption, for example, they are not stochastically ordered, the PO model may be misleading.

3.2. Non-Proportional Odds Model

To relax the assumption required by the PO model, that the log-odds ratios do not differ with the response level, we propose a hierarchical cumulative logistic regression model that assumes

$$\text{logit} \{ \pi_y^+(X, A; \alpha_y, \gamma_y) \} = \alpha_y + \gamma_{1,y} X + \gamma_{2,y} A + \gamma_{3,y} XA,$$

for $y = 0, \dots, 4$, (4)

where X and A are defined similarly as for the PO model. In contrast with the PO model, the NPO model in (4) allows different log-odds ratios at each response level y . We assume that observations are mutually independent such that the likelihood function for the unknown parameters (α, γ) has the same general form (2) as for the PO model.

We specify a hierarchical prior for (α, γ) as follows,

$$\alpha_0 \sim t_5(\alpha_0^*, 2.5), \alpha_y | \alpha_{y-1}, \gamma_{y-1}, \gamma_y \sim t_5(\alpha_y^*, 2.5) [\alpha_{y-1} + 0.5 | \gamma_{1,y-1} - \gamma_{1,y} | + 0.5 | \gamma_{2,y-1} - \gamma_{2,y} | + 0.25 | \gamma_{3,y-1} - \gamma_{3,y} |, \infty],$$

for $y = 1, 2, 3, 4$,

$$\gamma_{1,y} | \beta_1, \sigma_1 \sim N(\beta_1, \sigma_1^2), \beta_1 \sim t_5(\beta_1^*, 2.5), \sigma_1 \sim N(0, 1) [0, \infty],$$

$$\gamma_{2,y} | \beta_2, \sigma_2 \sim N(\beta_2, \sigma_2^2), \beta_2 \sim t_5(0, 2.5), \sigma_2 \sim N(0, 1) [0, \infty],$$

$$\gamma_{3,y} | \beta_3, \sigma_3 \sim N(\beta_3, \sigma_3^2), \beta_3 \sim t_5(0, 2.5), \sigma_3 \sim N(0, 1) [0, \infty].$$

(5)

The prior constraints on $\{\alpha_y\}_{y=1}^4$ ensure that $0 \leq \pi_1^+(X, A; \alpha_1, \gamma_1) \leq \dots \leq \pi_4^+(X, A; \alpha_4, \gamma_4) \leq 1$ for all $X \in \{-0.5, 0.5\}$ and $A \in \{-0.5, 0.5\}$. These inequalities hold when

$$\alpha_y \geq \alpha_{y-1} + 0.5 | \gamma_{1,y-1} - \gamma_{1,y} | + 0.5 | \gamma_{2,y-1} - \gamma_{2,y} | + 0.25 | \gamma_{3,y-1} - \gamma_{3,y} |, \quad y = 1, 2, 3, 4,$$

which is reflected in our specification of the prior. The hierarchical structure that we propose in (5) shrinks toward the PO model defined in (1). As $\sigma_1^2 \rightarrow 0$, $\sigma_2^2 \rightarrow 0$, and $\sigma_3^2 \rightarrow 0$, then $\gamma_{1,y} \rightarrow \beta_1$, $\gamma_{2,y} \rightarrow \beta_2$, and $\gamma_{3,y} \rightarrow \beta_3$, for $y = 0, \dots, 4$, and thus the log-odds ratios become invariant to the outcome level. Essentially, the NPO model in (4) adds a layer of additional structure to the PO model in (1) that allows each effect to deviate from the PO assumption. By using half-normal distributions with a unit standard deviation for σ_1 , σ_2 , and σ_3 ,

a priori the proposed NPO model prefers small deviations from the PO model. This type of hierarchical NPO model was alluded to as an alternative for PO models in the discussion of McKinley et al. (2015), but they neither implemented nor fully specified such a model. Because the proposed NPO model allows the log-odds ratios to differ with the response level, the model-based estimates of the response distributions need not be stochastically ordered. Consequently, when posterior inference is based on the NPO model, our proposed utility-based comparative testing criterion facilitates treatment comparison.

3.3. Posterior Estimation

We carry out posterior estimation for the PO and NPO models using JAGS via the R package R2jags (Plummer, 2003). Posterior convergence tends to be immediate and autocorrelation tends to be low, likely due, in part, to our balanced specification of the design matrix. We use the posterior samples from JAGS to calculate the four posterior probabilities required for the utility-based comparative testing criterion given in Section 2. Because the mean utilities are tractable functions of the unknown model parameters for both the PO and NPO models, obtaining posterior samples of the mean utilities is straightforward. We provide freely-available, user-friendly R software for implementation, see Supplementary Materials.

4. Design Considerations

Although we use a Bayesian probability model for statistical inference, we design the trial to ensure certain desirable frequentist operating characteristics (OCs), for example, 0.80 power under the targeted alternative with 0.05 probability of making a type I error. We are concerned with the subgroup specific power and type I error probability. That is, when Nuprehab reduces the number and severity of postoperative complications for a particular subgroup by the targeted amount, we want our design to have 0.80 probability of correctly declaring N superior to C for patients in that subgroup. By contrast, when Nuprehab does not reduce the number and severity of postoperative complications for a particular subgroup, we want our design only to have 0.05 probability of incorrectly declaring N superior or inferior to C for patients in that subgroup.

We met with WH to determine a practical targeted difference between the treatment arms for primary and salvage patients. During our discussion, WH expressed his desire for the trial to be powered to detect a 75% reduction in POM scores ≥ 3 in each subgroup. We then derived a targeted mean utility difference corresponding to this 75% reduction in POM scores ≥ 3 in each subgroup as follows. For the anticipated POM score probabilities in the control arm for each subgroup, a 75% reduction corresponds to shifts in the probability of POM scores ≥ 3 from 0.20 to 0.05 for primary patients, and from 0.35 to 0.09 for salvage patients. Under the proportional odds assumption, a 75% reduction in POM scores ≥ 3 corresponds to the cumulative POM score probabilities and mean utilities in Table 2. We designed the trial to target mean utility differences of $17.3 = 92.8 - 75.5$ in primary patients and $27.6 = 87.6 - 60.0$ in salvage patients.

Whitehead (1993) derived a sample size formula for a traditional fixed-sample design with ordinal outcomes based on a PO model, which is

$$n = 12 \left[\Phi^{-1}(1 - \alpha/2) + \Phi^{-1}(1 - \beta) \right]^2 / \left\{ \delta^2 \left[1 - \sum_{y=0}^5 (\pi_y^*)^3 \right] \right\},$$

where $\Phi(x)$ is the standard normal distribution evaluated at x , α , and β are the desired type I and II error probabilities, δ is the targeted log-odds ratio, and $\pi_y^* = \Pr(Y = y)$ under the targeted alternative. Although we use our utility-based comparative testing criterion proposed in Section 2, we demonstrated earlier that for the PO model this criterion is equivalent to a particular contrast of the regression coefficients that is also the basis for Whitehead's sample size formula. Viewing each subgroup as a separate trial, under their respective targeted alternatives, we need to enroll 56 primary and 38 salvage patients for $\alpha = 0.05$ and $\beta = 0.20$. Assuming 60% of the enrollees will be primary patients, we need to enroll 94 patients to achieve these subgroup sample sizes. Because the PO model borrows strength across subgroups for estimating the intercept parameters, $\{\alpha_y\}_{y=0}^4$, we expect the above sample size calculation to be conservative. When the PO assumption holds, we expect our NPO model to have less power than our PO model, though only slightly less as we specify an informative half-standard normal distribution as the prior for σ_1 , σ_2 and σ_3 in (5).

In the trial, patients are assigned to the two treatment arms using stratified block randomization with blocks of size four. Thus, for each block of four patients within each subgroup, two patients will receive Nuprehab and two will receive control. This will ensure that the treatment arms will have similar numbers of patients from each subgroup throughout the trial. Given the modest sample size requirements, one interim analysis will be done half-way through the trial. At this point, using our utility-based comparative testing criterion, the design will decide whether to continue enrolling patients from each subgroup. If one treatment is declared superior to the other for a certain subgroup at the interim analysis, then no additional patients from that subgroup will be enrolled. Otherwise, enrollment of patients from that subgroup will continue until the final analysis.

To control the probability of committing a type I error, we use a maximum duration alpha-spending approach such that $f(t) = \alpha \times (t/T_{max})^3$, where T_{max} denotes the maximum trial duration, see Jennison and Turnbull (1999, Section 7.2.3). To do this, we set $p_{cur} = \Phi(z)$ where z corresponds to the relevant threshold for the test statistic in a frequentist group sequential analysis, which we calculate using the R package `gsDesign`. With one interim analysis at the mid-point of the trial, this gives probability thresholds at the interim and final analyses of 0.997 and 0.976, respectively. Due to the asymptotic normality of the posterior distribution in general, see, for example, Gelman et al. (2014, Section 4), these thresholds control the type I error asymptotically. However, we use computer simulation to verify that our design controls type I error for the planned sample size. To be conservative, we aim to enroll up to 100 patients. Given the anticipated accrual rate of two patients per month, the interim and final analyses are expected to be performed at 26 and 51 months, respectively.

Table 2

Targeted alternative for each subgroup in terms of the cumulative POM score probabilities and the corresponding mean utility

Subgroup	Treatment arm	POM score					Mean utility
		≤0	≤1	≤2	≤3	≤4	
Primary	Nutritional prehabilitation	0.83	0.92	0.95	0.98	0.99	92.8
Primary	Control	0.50	0.70	0.80	0.90	0.95	75.5
Salvage	Nutritional prehabilitation	0.71	0.87	0.91	0.94	0.97	87.6
Salvage	Control	0.30	0.55	0.65	0.75	0.85	60.0

5. Simulation Study

In this section, we describe a simulation study that we carried out to evaluate and compare the frequentist OCs of the proposed design, under each of the PO and NPO models defined in Section 3. For further comparison, we considered a one-size-fits-all design based on a PO model that assumes

$$\text{logit} \{ \pi_y^+(X, A; \alpha_y, \beta) \} = \alpha_y + \beta_1 X + \beta_2 A, \text{ for } y = 0, \dots, 4.$$

We specify the same prior distributions for α , β_1 , and β_2 as for the PO model defined in (3). For this design, using the same two-stage group sequential structure that controls the probability of committing a type I error at 0.05, if $\text{Prob}(\beta_2 > 0 | \mathbf{D}) > p_{cut}$, then we declare N superior to C for patients in both subgroups. Conversely, if $\text{Prob}(\beta_2 < 0 | \mathbf{D}) > p_{cut}$, then we declare N inferior to C for patients in both subgroups. We compare the designs based on their probabilities of declaring N superior (or inferior) to C across a range of scenarios. To assess the decision criteria, we used 10,000 posterior samples following 500 warm-up samples. For all three models, posterior sampling took about 4 seconds for an interim analysis

with 50 observations, and about 7 seconds for a final analysis with up to 100 observations.

Table 3 reports the true POM score distributions for each scenario. We used the same POM score distribution to generate observations for the control arm in each subgroup for every scenario, that is, \mathbf{P}_0 for primary patients and \mathbf{S}_0 for salvage patients. Each patient had a 60% chance of belonging to the primary subgroup throughout. Scenario 1 is the complete null case where N provides no benefit to patients in either subgroup. Scenarios 2 and 3 are treatment-subgroup interaction cases where N provides the targeted benefit to patients in one subgroup, and no benefit to patients in the other subgroup. Scenario 4 is the complete alternative case where N provides the targeted benefit to patients in both subgroups. The PO assumption holds for Scenarios 1–4. Scenarios 5 and 6 are treatment-subgroup interaction cases, and Scenario 7 is a complete alternative case, but the PO assumption does not hold. In particular, the log-odds ratio comparing treatment arms corresponding to POM scores ≤ 0 and ≤ 1 are smaller than those corresponding to POM scores ≤ 2 , ≤ 3 , and ≤ 4 , but the targeted 75% reduction in POM scores ≥ 3 is still achieved in each subgroup. This reflects a benefit that greatly

Table 3

Simulation scenarios defined by the POM score distributions in each subgroup and treatment arm, and the corresponding mean utility

True POM Score distribution	POM score					Mean utility
	≤0	≤1	≤2	≤3	≤4	
\mathbf{P}_0	0.50	0.70	0.80	0.90	0.95	75.5
\mathbf{P}_1	0.83	0.92	0.95	0.98	0.99	92.8
\mathbf{P}_2	0.67	0.85	0.95	0.98	0.99	88.8
\mathbf{S}_0	0.30	0.55	0.65	0.75	0.85	60.0
\mathbf{S}_1	0.71	0.87	0.91	0.94	0.97	87.6
\mathbf{S}_2	0.53	0.80	0.91	0.94	0.97	82.8

Scenario	POM score distribution				Mean utility difference	
	P,C	P,N	S,C	S,N	P	S
1	\mathbf{P}_0	\mathbf{P}_0	\mathbf{S}_0	\mathbf{S}_0	0.0	0.0
2	\mathbf{P}_0	\mathbf{P}_1	\mathbf{S}_0	\mathbf{S}_0	17.3	0.0
3	\mathbf{P}_0	\mathbf{P}_0	\mathbf{S}_0	\mathbf{S}_1	0.0	27.6
4	\mathbf{P}_0	\mathbf{P}_1	\mathbf{S}_0	\mathbf{S}_1	17.3	27.6
5	\mathbf{P}_0	\mathbf{P}_2	\mathbf{S}_0	\mathbf{S}_0	13.3	0.0
6	\mathbf{P}_0	\mathbf{P}_0	\mathbf{S}_0	\mathbf{S}_2	0.0	22.8
7	\mathbf{P}_0	\mathbf{P}_2	\mathbf{S}_0	\mathbf{S}_2	13.3	22.8

Table 4

Simulation results. We report the proportion of simulated trials in which N was declared superior (inferior) to C , and the average sample size. The correct decision is indicated in boldface. Reported figures are based on 5000 simulated trials.

Design	Primary		Salvage		Average sample size
	Sup.	Inf.	Sup.	Inf.	
Scenario 1 (Complete null)					
Traditional	0.022	0.025	0.022	0.025	99.7
Stratified PO	0.022	0.025	0.020	0.025	99.8
Stratified NPO	0.029	0.028	0.029	0.032	99.5
Scenario 2 (Partial null-PO)					
Traditional	0.509	0.000	0.509	0.000	95.5
Stratified PO	0.774	0.000	0.038	0.017	94.9
Stratified NPO	0.763	0.000	0.048	0.022	93.7
Scenario 3 (Partial null-PO)					
Traditional	0.417	0.000	0.417	0.000	96.6
Stratified PO	0.040	0.015	0.784	0.000	96.5
Stratified NPO	0.047	0.016	0.782	0.000	95.4
Scenario 4 (Complete alternative-PO)					
Traditional	0.978	0.000	0.978	0.000	74.7
Stratified PO	0.841	0.000	0.850	0.000	87.4
Stratified NPO	0.842	0.000	0.853	0.000	84.6
Scenario 5 (Partial null-NPO)					
Traditional	0.215	0.001	0.215	0.001	98.5
Stratified PO	0.314	0.001	0.032	0.019	98.9
Stratified NPO	0.347	0.001	0.043	0.025	98.0
Scenario 6 (Partial null-NPO)					
Traditional	0.239	0.001	0.239	0.001	98.3
Stratified PO	0.036	0.019	0.462	0.000	98.7
Stratified NPO	0.042	0.021	0.503	0.000	97.7
Scenario 7 (Complete alternative-NPO)					
Traditional	0.710	0.000	0.710	0.000	92.1
Stratified PO	0.387	0.000	0.528	0.000	96.6
Stratified NPO	0.434	0.000	0.584	0.000	94.8

reduces severe postoperative complications, but affects the rate of minor postoperative complications to a lesser degree.

Table 4 reports the proportion of trials in which N was declared superior (inferior) to C , and the average sample size. In Scenario 1, that is, the complete null case, all three designs control the within subgroup probability of making a type I error near 0.05. In Scenario 4, that is, the complete alternative case, the two stratified designs have greater than 0.80 power of declaring N superior to C in each subgroup. Excepting Scenario 4, the three competing designs are unlikely to stop early, which is reflected by the average sample sizes near 100. Compared to the stratified design based on the PO model, when the PO assumption holds, for example, Scenarios 2, 3, and 4, the stratified design based on the more flexible NPO model

has similar power for declaring N superior to C in the benefiting subgroup(s), and when the PO assumption does not hold, for example, Scenarios 5, 6, and 7, the NPO model has larger power of declaring N superior to C in the benefiting subgroup(s). The power for each design is lower when the PO assumption does not hold, which is not surprising as the mean utility differences are smaller in these cases. Compared to the traditional design, when both subgroups benefit, for example, Scenarios 4 and 7, the stratified designs have less power for declaring N superior to C in both subgroups, but when only one subgroup benefits, for example, Scenarios 2, 3, 5, and 6; the stratified designs have greater power of declaring N superior to C in the benefiting subgroup, and are far less likely to incorrectly declare N superior to C in the non-benefiting subgroup.

6. Discussion

In this article, we have proposed a design for comparing treatments in two prognostic subgroups based on ordinal outcomes. The design was motivated by a trial comparing the effectiveness of nutritional prehabilitation (NuPrehab) against the standard of care for improving postoperative outcomes in primary and salvage patients who undergo esophageal resection. We considered two Bayesian cumulative logistic regression models for statistical inference, a proportional odds (PO) model and a hierarchical non-proportional odds (NPO) model that shrinks toward the PO model. Based on the results of our simulation study, we determined that the design based on the NPO model has preferable frequentist operating characteristics. In particular, when the PO assumption is satisfied the NPO model has similar probability of recommending NuPrehab in the benefiting subgroup(s), whereas when the PO assumption is not satisfied the NPO model can have substantially higher probability of recommending NuPrehab in the benefiting subgroup(s).

We also compared our stratified medicine design with a more traditional design that does not facilitate subgroup specific recommendations. If both subgroups benefit from NuPrehab, then the traditional design is more likely to recommend its use in each subgroup. However, when only one of the two subgroups benefits from NuPrehab, the traditional design is less likely to recommend its use in the benefiting subgroup and more likely to recommend its use in the non-benefiting subgroup. Because substantial treatment effect heterogeneity between primary and salvage patients is plausible in the motivating trial, we find the stratified medicine design based on the NPO model appealing. In another context, if this design is not appealing, then including a model selection between a model with and a model without an interaction term may provide an avenue for achieving a design with more appealing operating characteristics. Another alternative is to consider a prior for the interaction parameter(s) that facilitates borrowing substantial strength for the treatment effect across subgroups, perhaps in a data-dependent manner. These are avenues for future research.

7. Supplementary Materials

Web Appendices with a proof of Theorem 1 referenced in Section 2, and R software referenced in Section 3 for implementing the probability models and reproducing the simulation study is available with this article at the *Biometrics* website on Wiley Online Library.

ACKNOWLEDGEMENT

Peter F. Thall was funded in part by NIH/NCI Grant 5-R01-CA083932

REFERENCES

Abreu, M. N. S., Siqueira, A. L., Cardoso, C. S., and Caiaffa, W. T. (2008). Ordinal logistic regression models: Application in quality of life studies. *Cadernos de Saúde Pública* **24**, 581–591.

Agresti, A. (2010). *Analysis of Ordinal Categorical Data*. Hoboken, NJ: John Wiley & Sons, Inc.

Bender, R. and Grouven, U. (1998). Using binary logistic regression models for ordinal data with non-proportional odds. *Journal of Clinical Epidemiology* **51**, 809–816.

Braga, M., Gianotti, L., Nespoli, L., Radaelli, G., and Di Carlo, V. (2002). Nutritional approach in malnourished surgical patients: A prospective randomized study. *Archives of Surgery* **137**, 174–180.

Brannath, W., Zuber, E., Branson, M., Bretz, F., Gallo, P., Posch, M., et al. (2009). Confirmatory adaptive designs with bayesian decision tools for a targeted therapy in oncology. *Statistics in Medicine* **28**, 1445–1463.

Chen, K.-N. (2014). Managing complications I: Leaks, strictures, emptying, reflux, chylothorax. *Journal of Thoracic Disease* **6**, 355–363.

Clavien, P. A., Barkun, J., de Oliveira, M. L., Vauthey, J. N., Dindo, D., Schulick, R. D., et al. (2009). The Clavien-Dindo classification of surgical complications. *Annals of Surgery* **250**, 187–196.

Clavien, P. A., Sanabria, J. R., and Strasberg, S. M. (1992). Proposed classification of complications of surgery with examples of utility in cholecystectomy. *Surgery* **111**, 518–526.

Dindo, D., Demartines, N., and Clavien, P.-A. (2004). Classification of surgical complications: A new proposal with evaluation in a cohort of 6336 patients and results of a survey. *Annals of Surgery* **240**, 205–213.

Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., and Rubin, D. B. (2014). *Bayesian Data Analysis, 3rd edition*. Boca Raton, FL: Chapman & Hall/CRC Press.

Ghosh, J., Li, Y., and Mitra, R. (2017). On the use of Cauchy prior distributions for Bayesian logistic regression. *Bayesian Analysis* 1–25.

Guo, B. and Yuan, Y. (2017). Bayesian phase I/II biomarker-based dose finding for precision medicine with molecularly targeted agents. *Journal of the American Statistical Association* **112**, 508–520.

Houede, N., Thall, P. F., Nguyen, H., Paoletti, X., and Kramar, A. (2010). Utility-based optimization of combination therapy using ordinal toxicity and efficacy in phase I/II trials. *Biometrics* **66**, 532–540.

Ishwaran, H. (2000). Univariate and multivariate ordinal cumulative link regression with covariate specific cutpoints. *Canadian Journal of Statistics* **28**, 715–730.

Jennison, C. and Turnbull, B. W. (1999). *Group Sequential Methods with Applications to Clinical Trials*. Boca Raton, FL: Chapman and Hall/CRC.

McCullagh, P. (1980). Regression models for ordinal data. *Journal of the Royal Statistical Society, Series B (Methodological)* **42**, 109–142.

McKinley, T. J., Morters, M., and Wood, J. L. N. (2015). Bayesian model choice in cumulative link ordinal regression models. *Bayesian Analysis* **10**, 1–30.

Murray, T. A., Thall, P. F., and Yuan, Y. (2016). Utility-based designs for randomized comparative trials with categorical outcomes. *Statistics in Medicine* **35**, 4285–4305.

Parekh, K. and Iannettoni, M. D. (2007). Complications of esophageal resection and reconstruction. *Seminars in Thoracic and Cardiovascular Surgery* **19**, 79–88.

Peterson, B. and Harrell, F. E. (1990). Partial proportional odds models for ordinal response variables. *Journal of the Royal Statistical Society, Series C* **39**, 205–217.

Plummer, M. (2003). JAGS: A program for analysis of Bayesian graphical models using Gibbs sampling.

Quirk, J. P. and Saposnik, R. (1962). Admissibility and measurable utility functions. *The Review of Economic Studies* **29**, 140–146.

- Rosenblum, M., Luber, B., Thompson, R. E., and Hanley, D. (2016). Group sequential designs with prospectively planned rules for subpopulation enrichment. *Statistics in Medicine* **35**, 3776–3791. sim.6957.
- Thall, P. F. and Nguyen, H. Q. (2012). Adaptive randomization to improve utility-based dose-finding with bivariate ordinal outcomes. *Journal of Biopharmaceutical Statistics* **22**, 785–801.
- Waitzberg, D. L., Saito, H., Plank, L. D., Jamieson, G. G., Jagannath, P., Hwang, T.-L., et al. (2006). Postsurgical infections are reduced with specialized nutrition support. *World Journal of Surgery* **30**, 1592–1604.
- Walters, S., Campbell, M., and Lall, R. (2001). Design and analysis of trials with quality of life as an outcome: A practical guide. *Journal of Biopharmaceutical Statistics* **11**, 155–176.
- Wang, S.-J., James Hung, H. M., and O’Neill, R. T. (2009). Adaptive patient enrichment designs in therapeutic trials. *Biometrical Journal* **51**, 358–374.
- Whitehead, J. (1993). Sample size calculations for ordered categorical data. *Statistics in Medicine* **12**, 2257–2271.

Received June 2017. Revised November 2017.

Accepted November 2017.