# SIS | 2023

## Società Italiana di Statistica

**Statistical LEArning, Sustainability and Impact EvaluatioN**

## SEAS IN

Università Politecnica delle Marche

June 21-23, 2023
Ancona (Italy)

# Book of the Short Papers

**Editors: Francesco Maria Chelli, Mariateresa Ciommi, Salvatore Ingrassia, Francesca Mariani, Maria Cristina Recchioni**

SIS Società Italiana di Statistica

UNIVERSITÀ POLITECNICA DELLE MARCHE

Dipartimento di Scienze Economiche e Sociali DISES

LIUC Università Cattaneo

LIUC Università Cattaneo | BUSINESS ANALYTICS AND DATA SCIENCE HUB

Pearson

CHAIRS

Salvatore Ingrassia (Chair of the Program Committee) - *Università degli Studi di Catania*
Maria Cristina Recchioni (Chair of the Local Organizing Committee) - *Università Politecnica delle Marche*

PROGRAM COMMITTEE

Salvatore Ingrassia (Chair), Elena Ambrosetti, Antonio Balzanella, Matilde Bini, Annalisa Busetta, Fabio Centofanti, Francesco M. Chelli, Simone Di Zio, Sabrina Giordano, Rosaria Ignaccolo, Filomena Maggino, Stefania Mignani, Lucia Paci, Monica Palma, Emilia Rocco.

LOCAL ORGANIZING COMMITTEE

Maria Cristina Recchioni (Chair), Chiara Capogrossi, Mariateresa Ciommi, Barbara Ermini, Chiara Gigliarano, Riccardo Lucchetti, Francesca Mariani, Gloria Polinesi, Giuseppe Ricciardo Lamonica, Barbara Zagaglia.

ORGANIZERS OF INVITED SESSIONS

Pierfrancesco Alaimo Di Loro, Laura Anderlucci, Luigi Augugliaro, Ilaria Benedetti, Rossella Berni, Mario Bolzan, Silvia Cagnone, Michela Cameletti, Federico Camerlenghi, Gabriella Campolo, Christian Capezza, Carlo Cavicchia, Mariateresa Ciommi, Guido Consonni, Giuseppe Ricciardo Lamonica, Regina Liu, Daniela Marella, Francesca Mariani, Matteo Mazziotta, Stefano Mazzuco, Raya Muttarak, Livia Elisa Ortensi, Edoardo Otranto, Ilaria Prosdocimi, Pasquale Sarnacchiaro, Manuela Stranges, Claudia Tarantola, Isabella Sulis, Roberta Varriale, Rosanna Verde.

FURTHER PEPOLE OF LOCAL ORGANIZING COMMITTEE

Elisa D'Adamo, Christian Ferretti, Giada Gabbianelli, Elvina Merkaj, Luca Pedini, Alessandro Pionati, Marco Tedeschi, Francesco Valentini, Rostand Arland Yebetchou Tchounkeu

Technical support: Matteo Mercuri, Maila Ragni, Daniele Ripanti

# Contents

III

## Statistical Machine Learning for environmental applications

## Statistical Process Monitoring for Complex Data in Industry 4.0

## Advances in Data Science and Statistical Learning [IMS Invited Session]

## ENBIS Session: System Maintenance, Boosting algorithms for regression, and Research Excellence

# Population Dynamics, Climate Change and Sustainability

# Statistical Learning for health research and omics data

# The Economic behaviour of Sustainability

# Advances in statistical methods for complex problems

# Explainable machine learning models

## Flexible Learning for Environmental Sustainability

## Inequalities in higher education outcomes: learning from data

## Statistical Learning of demographic and health dynamics

## Challenges towards Fairness and Transparency for Data Processes, Algorithms and Decision-Support Models

# 3  Contributed Sessions  305

## Bayesian nonparametric methods

## Economics and Statistics

## Health statistics 1

## Indicators: composition, uses and limitations

## Multivariate data analysis 1

## Statistics in Society 1

## Assessment and Education

## Bayesian methods and applications 1

## Health and mortality

## Clustering and classification 1

## Dynamic models and time series

## Environmental learning and indicators

## Health statistics 2

## Migrants in Italy and return migration

## Sustainability assessment

## Bayesian methods and applications 2

## Clustering and classification 2

## Economics and labour markets

## Bayesian methods and applications 3

## Functional Data Analysis

## Machine Learning and text mining

## Multivariate data analysis 3

# Preface

This book includes the contributions presented at the Intermediate Meeting of the Italian Statistical Society (SIS) "SIS 2023 - Statistical Learning, Sustanaibility and Impact Evolution" held in Ancona at the Università Politecnica delle Marche, from June 21th to 23th of 2023.

The new challenges of digitalization, innovation and sustainability are showing the crucial role of data-driven approaches in supporting decision-making processes. Methodologies resulting from the integration of different know-how seem to be a reliable way to deal with the increasing need to measure the impact of the policies and to forecast scenarios. This meeting welcomed any attempt to face new challenges.

The conference registered more than 250 presentations, including 3 keynote speakers in 3 plenary sessions and 72 presentations in 24 invited sessions, all dealing with specific themes in methodological and/or applied statistics and demography. Furthermore, more than 180 contributions, with one or more authors, have been spontaneously submitted to the Program Committee and arranged in 30 contributed sessions.

The numerous participation of researchers in the conference shows how the challenges of sustainability, in its broadest sense, are of interest to both methodological and applied statistics.

With the publication of this book, we wish to offer to all members of the Italian Statistical Society, all international academics, researchers, Ph.D. students, and all interested practitioners, a good snapshot of the on-going research in the statistical and demographic fields.
We aim to provide all members of the Italian Statistical Society - as well as international academics, researchers, Ph.D. students, and interested practitioners - with a comprehensive overview of the ongoing research in the fields of statistics and demography.

We extend our heartfelt gratitude to all the contributors for submitting their works to the conference and to the researchers for their outstanding job in serving as referees and discussants with precision and timeliness.

A special appreciation goes to the Scientific and Organizational Committees for their tremendous efforts in managing all the organizational aspects, as well as to the Università Politecnica delle Marche and the Department of Economic and Social Science for making this event possible.

Finally, we wish to express our gratitude to the publisher Pearson Italia for all the support received.

# Optimal two-stage design based on error rates under a Bayesian perspective

Susanna Gentile[a] and Valeria Sambucini[a]

[a]Dipartimento di Scienze Statistiche, Sapienza Università di Roma;
susanna.gentile@uniroma1.it, valeria.sambucini@uniroma1.it

## Abstract

In phase II clinical trials, two-stage designs allowing early stopping for lack of efficacy are frequently used. We present a Bayesian two-stage design that ensures high posterior probabilities that the response rate of the experimental drug exceeds a desirable level, when the decision is to proceed with treatment evaluation. Moreover, the design exploits the distinction between *analysis* and *design* prior distributions, to control the predictive probability of Type I and II errors, while minimizing the expected sample size under the null hypothesis.

*Keywords:* analysis and design priors, Bayesian approach, error rates, two-stage design

## 1. Introduction

Single-arm two-stage designs are frequently used in phase II clinical trials with binary endpoints. The aim is to stop the study early for futility if the response rate of the experimental drug is not sufficiently high. According to the original scheme suggested by Simon (1989), at the first stage $n_1$ patients are enrolled and if $s_1 \leq r_1$, the trial terminates for lack of efficacy, where $s_1$ denotes the observed number of responders. Otherwise, the trial continues to the second stage and $n_2$ additional patients are treated. Then, if the observed number of responders $s$ out of the total sample $n = n_1 + n_2$ is not grater than $r$, the trial stops and the treatment is declared not effective. Otherwise, it is recommended for a more rigorous evaluation in a phase III trial.

The most popular two-stage designs are the optimal and the minimax designs developed by Simon (1989) under a frequentist framework, with many extensions and modifications presented in the literature. Several Bayesian two-stage designs have been also proposed (see, among others, Tan and Machin, 2002; Sambucini, 2008; Dong et al., 2012; Matano and Sambucini, 2016). In a recent work, Shi and Yin (2018) proposed a Bayesian enhancement two-stage (BET) design, that ($i$) guarantees a high posterior probability of the response rate exceeding the target of interest, when the observed number of responders reaches the minimum required level to continue with the experimentation and ($ii$) controls the length of the HPD interval for the response rate. In line with Shi and Yin (2018), we propose a Bayesian two-stage design that satisfies condition ($i$) and also yields the minimum expected sample size under the null hypothesis, while controlling the probability of Type I and Type II errors at a certain prespecified levels. To compute the error probabilities, we use a predictive approach by exploiting different kinds of prior distributions, which play different roles in the design of the trial. More specifically, we introduce an *analysis prior distribution* to express pre-experimental information and to construct posterior distributions. On the other hand, we elicit two different *design prior distributions* to represent suitable design scenarios to evaluate the Type I and Type II error rates.

The outline of the article is as follows. In Section 2, we formalize the Bayesian problem. In Section 3, the procedure to obtain the predictive probabilities of errors is described and the Bayesian strategy to find the optimal design is illustrated. Some numerical results are given in Sections 4 and, finally, Section 5 contains a brief discussion.

## 2. Preliminaries

Let $\theta$ be the unknown response rate of the experimental treatment and assume that the interest is focused on testing $H_0 : \theta \leq \theta^*$ vs $H_1 : \theta > \theta^*$. Thus, the treatment is considered sufficiently promising if $\theta$ exceeds the target value $\theta^*$. Moreover, let us denote by $S_1$ and $S$ the total number of responders at the end of the first and the second stage, respectively.

We introduce a beta prior density, $\pi^A(\theta) = \text{Beta}(\theta; \alpha^A, \beta^A)$, and use it to obtain the posterior distributions of $\theta$ at the end of both the stages. $\pi^A(\theta)$ is called the *analysis prior distribution*, because it is used to represent pre-experimental knowledge available at the analysis stage, that can be for instance gathered by previous studies or derived from expert opinions. We have that $S_1|\theta \sim \text{Bin}(n_1, \theta)$ and, from standard conjugate analysis, the first stage posterior distribution of $\theta$ is $\text{Beta}(\theta; \alpha^A + s_1, \beta^A + n_1 - s_1)$. Since $S$ includes the number of successes of the first stage, the posterior distribution for $\theta$ at the end of the second stage should be conditional on the event $S_1 > r_1$. However, it is possible to show that this condition does not affect the second stage posterior distribution (see Sambucini, 2008), that results to be $\text{Beta}(\theta; \alpha^A + s, \beta^A + n - s)$.

In each stage, we assume that the trial terminates if the posterior probability assigned to the alternative hypothesis is not sufficiently high. More specifically, at the end of the first stage the trial proceeds to the second one if

$$Pr(\theta > \theta^*|S_1 = s_1, n_1) > \lambda_1, \tag{1}$$

and, similarly, at the end of the second stage we claim the experimental treatment promising if

$$Pr(\theta > \theta^*|S = s, n) > \lambda_2, \tag{2}$$

where $\lambda_1$ and $\lambda_2$ are two desired probability thresholds, typically fixed so that $\lambda_2 > \lambda_1$.

## 3. The proposed two-stage design

For fixed values of $n_1$ and $n$, the posterior probabilities in (1) and (2) are increasing functions of $s_1$ and $s$, respectively. We use an algorithm that searches the optimal design by varying $n$ from $n^{min} = 10$ to $n^{max} = 100$. For each value of $n$, we consider $n_1$ in the range from $n_1^{min} = \max\{5, \frac{n}{3}\}$ to $n_1^{max} = n-1$. This choice for $n_1^{min}$ aims at avoiding that the sample size of the first stage may be relatively small compared to the total sample size. For each couple of $n$ and $n_1$, the two-stage boundaries are selected as

$$r_1 = \min \left\{ s_1 \in \{0, ..., n_1\} : Pr(\theta > \theta^*|S_1 = s_1, n_1) > \lambda_1 \right\} - 1 \tag{3}$$

$$r = \min \left\{ s \in \{r_1 + 1, ..., n\} : Pr(\theta > \theta^*|S = s, n) > \lambda_2 \right\} - 1 \tag{4}$$

In this way, we obtain a set of designs $(n_1, r_1, n, r)$ such that in each stage the posterior probability assigned to the alternative hypothesis is larger than the threshold of interest, when the observed number of responders exceeds the corresponding boundary, $r_1$ or $r$. Among these designs, we select the one that satisfies error probability constraints at the end of the second stage and minimizes the expected sample size under $H_0$.

## 3.1 Error probability computation

We consider the standard two types of errors that can occur in hypothesis testing: rejecting the null hypothesis when it is actually true (Type I) and failing to reject the null hypothesis, when the alternative is

true (Type II). When studying the operating characteristic of two-stage designs under a frequentist framework, the probabilities of these errors are evaluated by specifying a single value for $\theta$, suitably selected under $H_0$ or $H_1$ according to the type of error we are interested in. By adopting a Bayesian approach, we instead introduce two different prior distributions that model uncertainty on the single values of the parameter specified in the classical framework and add flexibility to the procedure. In the statistical literature, these priors are typically called *design prior distributions*, because they are used at the design stage of the study to describe a scenario of interest and to derive the prior predictive distributions of the data (see, Wang and Gelfand, 2002; Sahu and Smith, 2006; Brutti et al., 2008; Sambucini, 2008). This allows to obtain the probability model that generates the data, under the assumption that $\theta$ is highly likely to be in a certain subset of the parameter space.

More specifically, to compute the Type I error rate, we need to realize the conjecture that $H_0$ is true. Thus, we introduce a beta design prior distribution for $\theta$, $\pi_{H_0}^D(\theta) = \text{Beta}(\theta; \alpha_{H_0}^D, \beta_{H_0}^D)$, which has mode $\theta_0^D$ smaller than $\theta^*$ and assigns negligible probability to values of the parameter under the alternative hypothesis. By marginalizing the sampling distribution of the data over $\pi_{H_0}^D(\theta)$, we obtain the prior predictive distribution of $S_1$ and $S - S_1$, that are

$$m_{H_0}^D(s_1) = \text{Bin-Beta}(s_1;\, n, \alpha_{H_0}^D, \beta_{H_0}^D), \quad \forall\, s_1 = 0, ..., n_1,$$

$$m_{H_0}^D(s - s_1) = \text{Bin-Beta}(s - s_1;\, n - n_1, \alpha_{H_0}^D, \beta_{H_0}^D), \quad \forall\, s - s_1 = 0, ..., n - n_1.$$

Here, $\text{Bin-Beta}(\cdot\,;\, m, a, b)$ denotes the probability mass function of a Binomial-Beta distribution with parameters $m$, $a$ and $b$. Therefore, given the four values $(n_1, r_1, n, r)$, the predictive probability of a Type I error is provided by

$$Pr(\text{Type I error}) = \sum_{i=r_1+1}^{n_1} \sum_{j=r+1}^{n} \text{Bin-Beta}(i;\, n_1, \alpha_{H_0}^D, \beta_{H_0}^D)\text{Bin-Beta}(j - i;\, n - n_1, \alpha_{H_0}^D, \beta_{H_0}^D).$$

When the focus is on the Type II error, we elicit a beta design prior distribution for $\theta$, $\pi_{H_1}^D(\theta) = \text{Beta}(\theta; \alpha_{H_1}^D, \beta_{H_1}^D)$, used to realize the assumption that the alternative hypothesis is true. Its prior mode is $\theta_1^D > \theta^*$ and the prior probability assigned to values of the $\theta$ under the null hypothesis is negligible. Analogously to the previuos case, $\pi_{H_1}^D(\theta)$ is exploited to obtain the prior predictive distribution of $S_1$ and $S - S_1$, that are still Binomial-Beta. Therefore, the predictive probability of a Type II error is given by

$$Pr(\text{Type II error}) = 1 - \sum_{i=r_1+1}^{n_1} \sum_{j=r+1}^{n} \text{Bin-Beta}(i;\, n_1, \alpha_{H_1}^D, \beta_{H_1}^D)\text{Bin-Beta}(j - i;\, n - n_1, \alpha_{H_1}^D, \beta_{H_1}^D).$$

## 3.2 Optimal design strategy

As described before, the first step of the proposed strategy to find the optimal design consists in considering all the possible couples of $n$ and $n_1$ and identifying the corresponding boundaries $r_1$ and $r$ by exploiting the conditions (3) and (4) about the posterior probabilities that the alternative hypothesis is true. Among the set of $(n_1, r_1, n, r)$ identified using this procedure, we select the one that

1. satisfies the error constraints $Pr(\text{Type I error}) < \alpha$ and $Pr(\text{Type II error}) < \beta$, where $\alpha$ and $\beta$ are desired levels for the error rates;
2. minimizes the expected sample size under $H_0$, $E(N|H_0) = n_1 + (n - n_1)(1 - PET(H_0))$, where $PET(H_0)$ is the *probability of early termination*, that is computed for $\theta$ equal to the mode of the design prior $\pi_{H_0}^D(\theta)$ and is given by

$$PET(H_0) = \sum_{i=0}^{r_1} \binom{n_1}{i}(\theta_0^D)^i(1 - \theta_0^D)^{n_1-i}.$$

## 4. Numerical results

In this Section, we report some numerical results to illustrate the main features of the proposed design. We set $\theta^* = 0.4$ and consider a non-informative analysis prior distribution, $\pi^A(\theta) = \text{Beta}(\theta; 1, 1)$. To elicit the design prior distributions, we resort to a typical way of proceeding by expressing the hyperparameters in terms of prior mode and prior sample size. For instance, the hyperparameters of the design prior $\pi^D_{H_1}(\theta)$ can be fixed as

$$\alpha^D_{H_1} = n^D_{H_1}\theta^D_1 + 1 \qquad \text{and} \qquad \beta^D_{H_1} = n^D_{H_1}(1 - \theta^D_1) + 1,$$

where the prior sample size $n^D_{H_1}$ regulates the concentration of the prior around its mode. In particular, we set $\theta^D_1 = 0.6$ and select $n^D_{H_1}$ so that it is equal to 0.99 the prior probability assigned to the intervals $[0.55, 0.65]$, $[0.5, 0.7]$ and $[0.4, 0.8]$. Consequently, we obtain three design priors with $n^D_{H_1}$ equal to 1035, 255 and 60, respectively, that assign negligible probability to values of $\theta$ smaller than $\theta^*$. These distributions are represented in Figure 1 along with the design prior $\pi^D_{H_0}(\theta)$. This latter density has mode $\theta^D_0 = 0.35$ and is based on a prior sample size, $n^D_{H_0}$, equal to 909. On the right of the graph, a small Table shows the optimal designs that correspond to the diffrent design priors $\pi^D_{H_1}(\theta)$. As expected, the corresponding optimal sample sizes increase when $n^D_{H_1}$ decreases, as a consequence of the greater dispersion of the design distribution.



| $\theta^*$ | $n^D_{H_1}$ | $r_1/n_1$ | $r/n$ |
|---|---|---|---|
| 0.4 | 1035 | 7/16 | 16/33 |
| | 255 | 7/16 | 19/40 |
| | 60 | 9/20 | 23/49 |

Figure 1: Design prior distributions $\pi^D_{H_0}(\theta)$ and $\pi^D_{H_1}(\theta)$ for different values $n^D_{H_1}$, when $\theta^* = 0.4$. The Table shows the corresponding optimal designs when $\pi^A(\theta) = \text{Beta}(\theta; 1, 1)$ and $(\alpha, \beta) = (0.05, 0.2)$

In Table 1, we report the optimal designs for different values of $\theta^*$ and $(\alpha, \beta)$, when $\theta^D_0 = \theta^* - 0.05$, $\theta^D_1 = \theta^* + 0.2$, $\alpha^A = \beta^A = 1$, $\lambda_1 = 0.8$ and $\lambda_2 = 0.9$. The prior sample sizes of the design priors, $n^D_{H_0}$ and $n^D_{H_1}$, are selected as the minimum values so that the probabilities assigned to $H_0$ and $H_1$, respectively, are equal to 0.999. The corresponding probabilities of early termination and the expected sample sizes are also reported. Regardless of $\theta^*$, the optimal sample sizes obtained when $\alpha = \beta = 0.1$ in both the stages are greater than the ones obtained with $\alpha = 0.05$ and $\beta = 0.1$. Furthermore, in the first case the probabilities of early termination are higher, resulting in expected sample sizes closer to $n_1$. As for the analysis prior distributions, we expect that their impact varies according to the degree of skepticism expressed towards the treatment efficacy.

Table 1: Optimal proposed two-stage design for different values of $\theta^*$, $\alpha$ and $\beta$, when $\theta_0^D = \theta^* - 0.05$, $\theta_1^D = \theta^* + 0.2$, $\alpha^A = \beta^A = 1$, $\lambda_1 = 0.8$ and $\lambda_2 = 0.9$

| $\theta^*$ | $(\alpha, \beta)$ | $r_1/n_1$ | $r/n$ | $EN(H_0)$ | $PET(H_0)$ |
|---|---|---|---|---|---|
| 0.2 | (0.1, 0.1) | 6/27 | 12/48 | 27.070 | 0.901 |
| | (0.05, 0.2) | 3/14 | 7/27 | 15.094 | 0.853 |
| 0.3 | (0.1, 0.1) | 11/33 | 25/70 | 36.652 | 0.901 |
| | (0.05, 0.2) | 6/18 | 14/38 | 20.780 | 0.861 |
| 0.4 | (0.1, 0.1) | 18/41 | 27/58 | 42.509 | 0.911 |
| | (0.05, 0.2) | 9/20 | 27/58 | 24.628 | 0.878 |
| 0.5 | (0.1, 0.1) | 22/40 | 40/71 | 42.377 | 0.923 |
| | (0.05, 0.2) | 11/20 | 27/47 | 23.531 | 0.869 |

## 5. Discussion

The most popular and commonly used two-stage designs have been developed by Simon (1989) under a frequentist framework. From a Bayesian perspective, Shi and Yin (2018) showed that, given promising results according to Simon's designs, the posterior probabilities that the response rate reaches the desirable target level are very low. Thus, this Author proposed a Bayesian two-stage design that ensures high posterior probabilities of the response rate exceeding the target of interest, when the observed results suggest to proceed with treatment investigation.

In this paper, we present a Bayesian two-stage design based on the same condition used by Shi and Yin (2018) and that in addition minimizes the expected sample size under null hypothesis, while controlling the Type I and Type II error rates at fixed desired levels. To offer a more flexible evaluation of the errors probabilities, we adopt a predictive approach by using design prior distributions to specify design expectations and to realize the assumption that $\theta$ belongs to a specific subset of the parameter space. These prior densities are employed to obtain the prior predictive distribution of the data, used to compute the error probabilities. A different prior distribution is used to represent prior information and to compute the posterior distribution of the parameter.

## References

[1] Brutti P., De Santis F., Gubbiotti S.: Robust Bayesian sample size determination in clinical trials. Stat. Med., **27**, 2290–2306 (2008)

[2] Dong G., Shih, W.J., Moore, D., Quand, H., Marcella, S.: A Bayesian-frequentist two-stage single-arm phase II clinical trial design. Stat. Med., **31**, 2055–2067 (2012)

[3] Matano, F., Sambucini, V.: Accounting for uncertainty in the historical response rate of the standard treatment in single-arm two-stage designs based on Bayesian power functions Pharm Stat., **15**, 517–530 (2016)

[4] Sahu, S.K., Smith, T.M.F.: A Bayesian method of sample size determination with practical applications. J. R. Stat. Soc., **169**, 235–253 (2006)

[5] Sambucini V.: A Bayesian predictive two-stage design for phase II clinical trials. Stat. Med., **27**, 1199–1224 (2008)

[6] Shih, H., Yin, G.: Bayesian enhancement two-stage design for single-arm phase II clinical trials with binary and time-to-event endpoints. Biometrics, **74**, 1055–1064 (2018)

[7] Tan, S. B., Machin, D.: Bayesian two stage designs for phase II clinical trials. Stat. Med., **21**, 1991–2012 (2002)

[8] Wang, F., Gelfand, A.E.: A simulation-based approach to Bayesian sample size determination for

performance under a given model and for separating models. Stat. Sci., **2**, 193–218 (2002)