# Code exploration

Minghan Yang

October 28, 2024

**Abstract**

Related to paper "Beyond p-values: A phase II dual-criterion design with statistical significance and clinical relevance" Roychoudhury et al. (2018).

## 1 Introduction

Proof-of-concept (POC) in Phase II trails is important in investigating the efficacy of an experimental drug. It will influence the decision of whether continuing or not continuing the development of this drug.

Dual-criterion design in frequentist and Bayesian applications are discussed.

Three generic phase II designs are reviewed:

1. Standard design

   For comparative treatment and control trails, it puts forward criteria expressed as error rates: Type I error control and Power (correctly reject $H_0$ when it is false).

   Control type I error and maximize power.

   Type I error: $\mathbb{P}(\text{reject } H_0 | H_0 \text{ is true}) = \alpha$

   Type II error: $\mathbb{P}(\text{not reject } H_0 | H_0 \text{ is false}) = \beta$.

   Limitation: statistical significant only guarantees evidence to reject "No effect", but is not sufficient for clinical perspective. Also, it always result in success or failure according to statistical significance.

   Increasing the sample size increases the power for effects better than null.

2. Dual-criterion design

   Considers both the statistical significance and the effect estimate.

   Required inputs: type I error control (null hypothesis and type I error $\alpha$) and a decision value (DV). The DV is same as the target difference in Fisch's paper. It is the minimal effect estimate needed for trail success (if higher than this value with moderate confidence, then GO).

   By considering both, we have both statistical significance and guarantees a sufficiently large effect estimate.

   The dual-criterion is more demanding, the resulting power of study is less compared to standard design.

   Power is only increased for values superior to the DV since inferior values are clinically irrelevant.

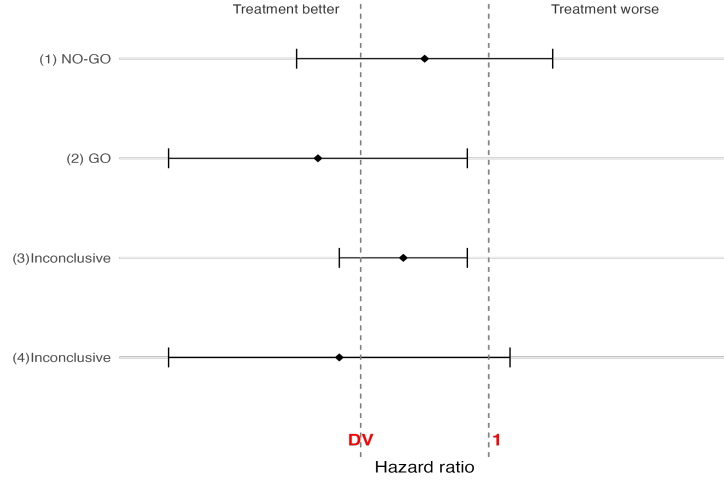   Decisions for dual-criterion design:

Figure 1: Decisions for dual criterion design

3. Precision design

Doesn't rely on error rates. When null hypothesis or other benchmark values cannot be determined, this can be an option. It requires sufficiently precise effect estimate.

Precision=$\frac{TP}{TP+FP}$. High Precision means that when the model predicts a positive outcome, it is very likely to be correct.

Consider hazard function in survival analysis, it describes the risk of failing. We consider hazard ratio between experimental drug and control as the outcome of interest. Hazard ratio(HR) less than 1 means the drug is better than control. We want to reduce hazard and hazard ratio.

**Sample size:**

$$n_{min} = \frac{\sigma^2 \times z_\alpha^2}{(NV - DV)^2}$$

It is the minimum sample size that implies statistical significance if the effect estimate equals the DV. This value is calculated under the situation that both criterion are just satisfied. As illustrated in the below graph, when the effect estimate $\theta = $ DV, and the lower bound of the confidence interval just touches the NV so that statistical significance is reached. Notice that when sample size equals the minimum sample size, there can only be GO or NO-GO decisions.
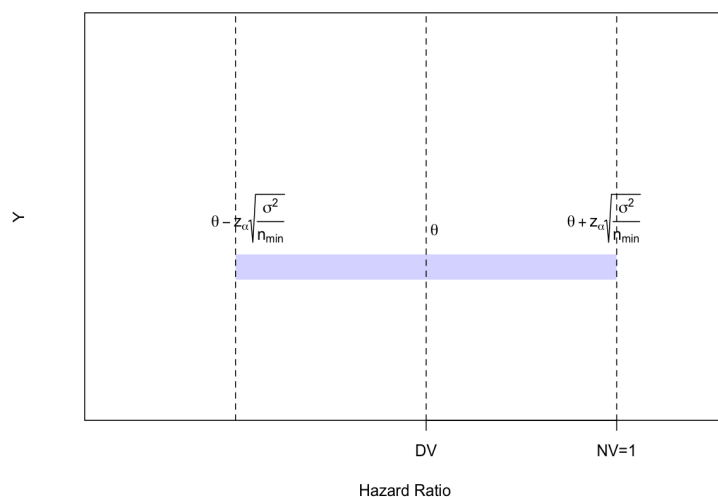
Figure 2: An illustration of sample size calculation

**Operating characteristics:**

The operating characteristics are the type I error and power of the clinical trial design.

For dual-criterion designs, the power at the DV is approximately 50%, so that if the true parameter equals the DV, there is roughly equal chance that the effect estimate lies on either side of the DV. Having 50% at the DV does not mean the study is under-powered.
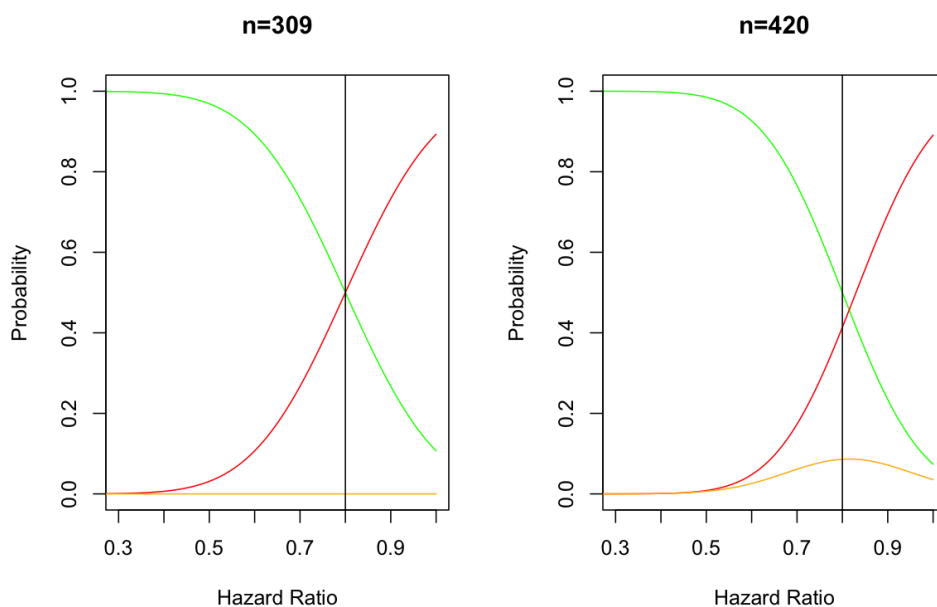
## 1.1 Reproduce Figure 1



Figure 3: Reproduce Figure 1

## 2 Example 1: A randomized PoC design with time-to-event data

Randomized, double-blind, RCT. Patients were randomized equally to: (experimental drug + standard care) OR (standard care only).

Primary outcome of interest (or called "endpoint") is the progression-free survival (PFS), which is the time when the disease or cancer do not get worse. The endpoint was assessed with a *log-rank test* and *Cox regression* with treatment as a covariate.

- *log-rank test*: compare the survival distributions of two or more groups. It tests the hypothesis that there is no difference in survival (or time-to-event) between the two groups. If the log-rank test indicates a significant difference, it suggests the treatment affects how long patients live without their disease worsening.

- *Cox regression*: estimate the hazard ratio between two groups, which tells us the relative risk of disease progression in the treatment group compared to the control group. If the $HR < 1$, it suggests that the new treatment delays disease progression better than the control treatment.

As for the DV, HR=0.7 was deemed necessary to be clinically meaningful. Values larger than 0.7 are unsatisfactory to clearly justify further development of the drug.

So the dual-criterion is:

1. Statistical significance: one-sided p-value of log-rank test $\leq 0.1$.

2. Clinical relevance: estimated HR from Cox regression $\leq 0.70$.

### 2.1 Reproduce Table 3

Attempt to reproduce the values in Table 3. The main problem is to find the correct value for $\sigma$. There seem to be no detailed information about the choice of value for $\sigma$. We chose it to be 2 according to page 455 in the paper. However, under this value, the calculated values do not agree with the ones in Table 3.

Also, the calculation of $\hat{\theta} = 0.736, 0.708, 0.659, 0.761$ should be looked into when the correct $\sigma$ is decided. The current understanding is that these values are calculated when the estimate $\hat{\theta} > DV$ and the confidence interval just touches 1, i.e. NO-GO is implied.

|   | True HR | GO | NO-GO | Inconclusive |
|---|---------|-----------|-----------|--------------|
| 1 | 0.5 | 0.7229434 | 0.2425581 | 0.03449843 |
| 1 | 0.6 | 0.6163101 | 0.3437340 | 0.03995584 |
| 1 | 0.7 | 0.5000000 | 0.4575970 | 0.04240298 |
| 1 | 0.8 | 0.3836899 | 0.5750769 | 0.04123326 |
| 1 | 0.9 | 0.2770566 | 0.6862039 | 0.03673956 |
| 1 | 1.0 | 0.1874286 | 0.7825760 | 0.02999546 |
| 2 | 0.5 | 0.6949399 | 0.3050601 | 0.00000000 |
| 2 | 0.6 | 0.6006195 | 0.3993805 | 0.00000000 |
| 2 | 0.7 | 0.5000000 | 0.5000000 | 0.00000000 |
| 2 | 0.8 | 0.3993805 | 0.6006195 | 0.00000000 |
| 2 | 0.9 | 0.3050601 | 0.6949399 | 0.00000000 |

| | | | | |
|---|---|---|---|---|
| 2 | 1.0 | 0.2221796 | 0.7778204 | 0.00000000 |
| 3 | 0.5 | 0.9000000 | 0.1000000 | 0.00000000 |
| 3 | 0.6 | 0.8459814 | 0.1540186 | 0.00000000 |
| 3 | 0.7 | 0.7755191 | 0.2244809 | 0.00000000 |
| 3 | 0.8 | 0.6896805 | 0.3103195 | 0.00000000 |
| 3 | 0.9 | 0.5920194 | 0.4079806 | 0.00000000 |
| 3 | 1.0 | 0.4882491 | 0.5117509 | 0.00000000 |
| 4 | 0.5 | 0.8000000 | 0.2000000 | 0.00000000 |
| 4 | 0.6 | 0.7335799 | 0.2664201 | 0.00000000 |
| 4 | 0.7 | 0.6575300 | 0.3424700 | 0.00000000 |
| 4 | 0.8 | 0.5744779 | 0.4255221 | 0.00000000 |
| 4 | 0.9 | 0.4879703 | 0.5120297 | 0.00000000 |
| 4 | 1.0 | 0.4020272 | 0.5979728 | 0.00000000 |
| 5 | 0.5 | 0.9000000 | 0.1000000 | 0.00000000 |
| 5 | 0.6 | 0.8562465 | 0.1437535 | 0.00000000 |
| 5 | 0.7 | 0.8011292 | 0.1988708 | 0.00000000 |
| 5 | 0.8 | 0.7349052 | 0.2650948 | 0.00000000 |
| 5 | 0.9 | 0.6590133 | 0.3409867 | 0.00000000 |
| 5 | 1.0 | 0.5760611 | 0.4239389 | 0.00000000 |

# 3 Example 2: A single-arm PoC design with binary data

Experimental drug in Chinese patients with non-small-cell lung cancer.

Primary endpoint is objective response rate (ORR), which quantifies the preliminary efficacy of the experimental drug.

Prior: minimally informative unimodal beta prior distribution $Beta(0.0811, 1)$, which has mean 0.75.

NV is set to 7.5% rather than 0, because of the absence of a comparator (in single arm trials).

DV is set to be 10%+7.5%=17.5%.

So the dual-criterion is:

1. Bayesian statistical significance: $\mathbb{P}(ORR \geq 7.5\%|data) \geq 0.95$

2. Clinical relevance: Posterior median $\geq 17.5\%$

The minimal sample size was 22. Final sample size 25.

Null hypothesis: there is no effect of the drug, i.e. ORR=7.5%

$\mathbb{P}(\text{type I error}) = \mathbb{P}(\text{reject } H_0|H_0 \text{ is true}) = \mathbb{P}(\text{reject } H_0|ORR \leq 7.5\%)$

$\mathbb{P}(\text{type II error}) = \mathbb{P}(\text{not reject } H_0|H_0 \text{ is false}) = \mathbb{P}(\text{reject } H_0|ORR = \text{response rate})$

Table 4 results show that this dual-criterion design is a three-outcome design with desirable properties.

## 3.1 Reproduce Figure 2

## 3.2 Reproduce Table 4

# References

Roychoudhury, S., Scheuer, N., , and Neuenschwander, B. (2018). Beyond p-values: A phase II dual-criterion design with statistical significance and clinical relevance. *Clinical trials (London, England)*, 15(5):452–461.

# Computational details

```
> cat(paste(Sys.time(), Sys.timezone(), "\n"))

2024-10-28 16:46:55.963281 Europe/Zurich

> sessionInfo()

R version 4.4.0 (2024-04-24)
Platform: aarch64-apple-darwin20
Running under: macOS Sonoma 14.4

Matrix products: default
BLAS:    /Library/Frameworks/R.framework/Versions/4.4-arm64/Resources/lib/libRblas.0.dylib
LAPACK: /Library/Frameworks/R.framework/Versions/4.4-arm64/Resources/lib/libRlapack.dylib;  LAI

locale:
[1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8

time zone: Europe/Zurich
tzcode source: internal

attached base packages:
[1] stats     graphics  grDevices utils     datasets  methods   base

other attached packages:
[1] latex2exp_0.9.6 ggplot2_3.5.1   knitr_1.47

loaded via a namespace (and not attached):
 [1] vctrs_0.6.5        cli_3.6.2          rlang_1.1.4        xfun_0.44
 [5] stringi_1.8.4      generics_0.1.3     textshaping_0.4.0  glue_1.7.0
 [9] colorspace_2.1-0   ragg_1.3.2         scales_1.3.0       fansi_1.0.6
[13] grid_4.4.0         munsell_0.5.1      tibble_3.2.1       lifecycle_1.0.4
[17] stringr_1.5.1      compiler_4.4.0     dplyr_1.1.4        pkgconfig_2.0.3
[21] rstudioapi_0.16.0  systemfonts_1.1.0  farver_2.1.2       R6_2.5.1
[25] tidyselect_1.2.1   utf8_1.2.4         pillar_1.9.0       magrittr_2.0.3
[29] tools_4.4.0        withr_3.0.0        gtable_0.3.5
```