# Study Design with Time to Event Endpoints

## Johnson and Johnson

**Overview**

Overview: There are ~25000 patient data across 5 indications (cancer type). The data includes clinical information (such as age, sex, treatment received, Overall survival etc.) and genetic data (e.g. mutation in a particular gene, what type of mutation in the gene, mutation location in the gene etc.)

**Data**

Download the data from cbioportal.org
(https://www.cbioportal.org/study/summary?id=msk_chord_2024)
(Download all clinical and genomic data)

**Background**

Often mutations in a gene in a cell surface led to disruption of apoptosis of cells. Our own body's immune system fails to recognize these abnormalities in cell life cycle which are often termed as malignant cells. A proliferation of these malignant cells is termed as a cancer of that organ where these cells are situated. Immunotherapy targets these mutations so that our immune system can again identify these abnormal cells and work to eliminate the malignant cells. There might be single or multiple mutations that can lead to different types of cancers. So, we try to target 1 or a combination of couple of mutations by an immunotherapy. Traditionally in most of the cases, a clinical trial is designed to find newer treatment to a particular type of cancer but if the same mutations can cause different type of cancers, we can perform a clinical trial with those different type of cancer (through different organs) driven by a particular mutation.

**Task 1**

The first problem is to identify a mutation that can be targeted and have sufficient patients in the database.

a.  Find out a gene mutation which is most prevalent in at least 2 out of 5 cancer types. Provide the methodology on how you have identified this genetic mutation target.
b.  Now subset the corresponding patients from the database and provide their baseline demographics and disease characteristics summary
      i.      Examples of demographic characteristics can be: Age, Sex, Ethnicity, Weight at baseline
      ii.     Example of disease characteristics can be: Number of prior therapies (prior anti-cancer medications, surgery etc.), summary of overall survival time etc.
c.  Comment on your observations.

**Task 2**

Our next problem will be based on designing a lung cancer trial. (The components of design includes setting a time-to-event endpoint such as Overall Survival, choosing proper hazard rates tostudy operating characteristics of the trial design and computing respective sample sizes etc.) We consider only subjects with lung cancer from the given dataset (Cancer.Type = "Non-Small Cell Lung Cancer") and will assume these subjects represent the population of our control group of the study.

a.  Using Kaplan-Meier calculate the median survival time. The columns to consider are "Overall.Survival..Months.", "Overall.Survival.Status".
b.  Assume exponential distribution for the survival time and establish the relation between hazard rate and median survival time. And compute the hazard rate for the population.
c.  Based on the hazard rate from (b) simulate n (size of the control group) survival times.
d.  For simplicity we will consider same size for treatment and control group. Simulate n (size of the treatment group) survival times based on the hazard rate from (b) and assuming a hazard ratio of 0.7.
e.  Using log-rank test (for the hypothesis hazard ratio < 1) compute the p-value.
f.  Repeat steps {c, d, e} 1000 times. The estimate of the power at is the proportion of times p-value $< 0.05$.
g.  Try a range of sample sizes (n in steps c and d) and comment on the minimum sample-size required to achieve 90% power. Show the relation between power and sample-size graphically.
h.  Thus far, we have used only hazard ratio of 0.7. In this step we will vary hazard ratio from 0.6 to 1. Comment on the minimum sample size required to achieve 90% power in each case.

**Task 3**

In this task we will enrich the study design obtained in Task-2, by introducing recruitment rate and dropout rate.

    a.  So far, the study design assumes all the participants were recruited simultaneously. However, in clinical trials patients are recruited over a span of years. Assume a constant recruitment rate (say 20 patients per month), and update the simulation set up. Note that you can assume uniform distribution for time of recruitment for 20 patients in a month. (Hint: time of death in calendar time is time of recruitment for a patient + time until death since recruitment of that patient).

    b.  We define total follow up time as the entire duration of time that participants are monitored and assessed for study outcomes. For the case where hazard ratio 0.7, choose a suitably large sample size (n) for each group, and comment on the minimum follow up time required to achieve a 90% power.

    c.  By varying hazard ratios comment on the follow up time required to achieve 90% power.

    d.  (Optional): We can also assume a dropout rate (subjects discontinuing the study without event). Assume a dropout rate of 1 person per month and update the simulation steps {a, b, c}.

       Note: In general, the budget of the study increases with the number of subjects. We are not considering budget in this exercise though.

**Task 4 (Optional)**

A trial may go on for a sufficiently long time. In many situations, we want to evaluate the primary research question of the trial in the *interim* with help of statistical review of the data collected till that time point. Often this is referred to as interim analysis. Assume that there is an interim analysis after 2 years and you are providing an estimate of the probability of response rate being more than 40% i.e., $P(\theta > 0.4)$, where $\theta$ denotes the response rate. (Responders are defined as when a patient has achieved Complete Response (CR) or Partial Response (PR); see RECIST paper attached.) This is different from overall survival data in Task 2 and is generally explored for internal decision-making purposes when we don't have full overall survival data or it's too long to wait for the overall survival data to mature. Assume $\theta \sim Bin(n, \pi)$.

The internal board has decided that the study will continue or to be stopped based on clinical benefit to the patients after 2 years. Clinical benefit is defined as $P(\theta > 0.4) > 0.8$. Assuming at least 3 different choices of response rate where,

         x1 := response rate of indication 1 (say, lung cancer),

x2 := response rate of indication 2 (say, colorectal cancer) and

x1 and x2 belong in [0.2, 0.8],

do the following:

a) Simulate response data based on sample size of n (Pick choices of sample size n1, n2 from task 2c) from Bin(n1, x1) and Bin(n2, x2)

b) Assume a weakly informative prior (e.g. Beta prior) for x1 and x2.

c) Find out the posterior probability estimate $P(\pi > 0.4|\text{data})$ for each scenario and indication. You may assume simple beta-binomial framework to calculate the posterior probability, or you can use hierarchical Bayesian modelling technique to find out the posterior probability. The choice is yours.

For hierarchical Bayesian framework, refer to Berry et al. BHM (primary), EXNEX Satrajit et al. (optional).

**References:**

1. Berry SM, Broglio KR, Groshen S, Berry DA. Bayesian hierarchical modeling of patient subpopulations: efficient designs of Phase II oncology clinical trials. Clin Trials. 2013 Oct;10(5):720-34. doi: 10.1177/1740774513497539. Epub 2013 Aug 27. PMID: 23983156; PMCID: PMC4319656.

2. Neuenschwander B, Wandel S, Roychoudhury S, Bailey S. Robust exchangeability designs for early phase clinical trials with multiple strata. Pharm Stat. 2016 Mar-Apr;15(2):123-34. doi: 10.1002/pst.1730. Epub 2015 Dec 18. PMID: 26685103.

3. Eisenhauer EA, Therasse P, Bogaerts J, Schwartz LH, Sargent D, Ford R, Dancey J, Arbuck S, Gwyther S, Mooney M, Rubinstein L, Shankar L, Dodd L, Kaplan R, Lacombe D, Verweij J. New response evaluation criteria in solid tumours: revised RECIST guideline (version 1.1). Eur J Cancer. 2009 Jan;45(2):228-47. doi: 10.1016/j.ejca.2008.10.026. PMID: 19097774.