8-2023

# Statistical inference on lung cancer screening using the national lung screening trial data.

Farhin Rahman
*University of Louisville*

# STATISTICAL INFERENCE ON LUNG CANCER SCREENING USING THE NATIONAL LUNG SCREENING TRIAL DATA

By

Farhin Rahman
B.S., University of Dhaka, 2011
M.S., Ball State University, 2017

A Dissertation Submitted to the Faculty of
the School of Public Health and Information Sciences
of the University of Louisville
in Partial Fulfillment of the Requirements
for the Degree of

Doctor of Philosophy
in Biostatistics

Department of Bioinformatics and Biostatistics
University of Louisville
Louisville, Kentucky

August, 2023

# STATISTICAL INFERENCE ON LUNG CANCER SCREENING USING THE NATIONAL LUNG SCREENING TRIAL DATA

By

Farhin Rahman
B.S., University of Dhaka, 2011
M.S., Ball State University, 2017

A Dissertation Approved on

July 19, 2023

By the following Dissertation Committee:

_____

Dongfeng Wu, Ph.D., Dissertation Director

_____

Jeremy Gaskins, Ph.D.

_____

Qi Zheng, Ph.D.

_____

Michael Sekula, Ph.D.

_____

Albert Seow, M.D.

# DEDICATION

I lovingly dedicate this dissertation to my dear parents, husband and my cherished daughter, whose boundless sacrifices have made my academic journey smooth and rewarding.

# ACKNOWLEDGMENTS

ABSTRACT


INFERENCE OF SOJOURN TIME AND TRANSITION
DENSITY USING THE NLST X-RAY SCREENING DATA IN
LUNG CANCER

Farhin Rahman

July 19, 2023

This dissertation consists of three research projects on cancer screening probability modeling. In these projects, the three key modeling parameters (sensitivity, sojourn time, transition density) for cancer screening were estimated, along with the long-term outcomes (including overdiagnosis as one outcome), the optimal screening time/age, the lead time distribution, and the probability of overdiagnosis at the future screening time were simulated to provide a statistical perspective on the effectiveness of cancer screening programs.

In the first part of this dissertation, a statistical inference was conducted for male and female smokers using the National Lung Screening Trial (NLST) chest X-ray data. A likelihood function was applied to the lung cancer screening data to obtain Bayesian inference of the transition probability and the distribution of sojourn time. For the transition probability density function, a log-normal distribution multiplied by 30% was used, while a Weibull distribution was employed to model the sojourn time in the preclinical state. The early transition occurred before age 50 and persisted until after age 90. Notably, the transition probability from the disease-free to the preclinical state peaked

around age 73 for males and 72 for females. The Bayesian posterior mean and median sojourn time for males (females) heavy smokers were estimated to be 1.28 (1.23) and 1.23 (1.21) years, respectively. These estimations revealed that male smokers are more susceptible to lung cancer due to their higher transition probability density compared to same-aged female smokers. Furthermore, female smokers exhibited a slightly shorter mean sojourn time than their male counterparts, suggesting they develop clinical symptoms of lung cancer at a faster rate.

In the second part of this dissertation, the probability model was applied to assess the long-term effects of cancer screening. The participants in the cancer screenings were categorized into four mutually exclusive groups: symptom-free-life, no-early-detection, true-early-detection, and overdiagnosis. Simulation studies and Bayesian inference were conducted, considering factors such as a person's age at the study entry, screening frequency, screening sensitivity, and other relevant parameters. The probability of overdiagnosis among the screen-detected cases was found to be relatively low but increased with the initial screen age. It was observed that males were more susceptible to overdiagnosis compared to females. The model can provide policymakers with essential information about the distribution of individuals in the overdiagnosis and true-early-detection groups, enabling them to minimize the long-term effects resulting from frequent screenings.

In the third part of this dissertation, a recently developed statistical method was applied to the National Lung Screening Trial (NLST) chest X-ray data, to find the optimal time for initiating chest X-ray screening in asymptomatic individuals. Incidence probability was used to control the risk of clinical incidence before the first exam, constraining it to a small value, given one's current age. The simulation study showed that the optimal screening

age interval remains relatively consistent as the current age increases. Notably, male heavy smokers tended to have slightly later screening ages compared to females, which contrasted with the findings from NLST CT data. Once the future screening time/age was found, the lead time distribution and the probability of overdiagnosis were estimated if one would be diagnosed at this future time/age. The lead time was relatively consistent across incidence probability and sensitivity, with a slight decrease in the mean lead time as the current age increased, and it was positively correlated with the sojourn time. The probability of overdiagnosis exhibited positive correlations with the mean sojourn time, incidence probability, and current age, and it only slightly changes with sensitivity. Overall, the probability of overdiagnosis was small and was not a concern at a younger age.

# TABLE OF CONTENTS

# LIST OF TABLES

# LIST OF FIGURES

# CHAPTER 1

# INTRODUCTION

This dissertation contains three research projects (Chapters 2-4) on cancer screening probability modeling. In Chapter 2, the Bayesian approach is employed to estimate the key parameters in cancer screening. Specifically, the sojourn time and transition density are estimated using the methodology described in a previous study by Wu et al. (2005). Chapter 3 focuses on estimating the probability of four outcomes, including overdiagnosis as one outcome in cancer screening. By incorporating factors such as age at study entry, screening frequency, and screening sensitivity, the probability of overdiagnosis is estimated using simulation studies and Bayesian inference techniques. In the third project, presented in Chapter 4, the focus shifts to finding the optimal screening time/age for an asymptomatic individual, then evaluating the lead time distribution and the probability of overdiagnosis at the future screening time. This is the first step toward individualized screening time for individuals at risk with helpful information for decision-making.

This chapter is a review of methods used to estimate the key parameters along with the optimal screening time/age, the lead time distribution, and the probability of overdiagnosis.

## 1.1 Lung Cancer Screening Overview

Lung cancer is a significant public health concern, being the second most common cancer and the leading cause of cancer-related deaths among both men and women in the United States (NCI, 2019). It arises when malignant (cancer) cells form in the lung tissues. The two main types of lung cancer are small cell and non-small cell, with the latter being more prevalent. These types exhibit different growth patterns and require distinct treatment approaches.

Cancer screening refers to the process of searching for cancer in individuals who do not exhibit any symptoms. The primary objective of lung cancer screening is to identify lung cancer in its early stages when it is more treatable and potentially curable. Early detection allows for timely intervention, such as surgical removal of tumors or other appropriate treatments, which can significantly improve outcomes and increase survival rates.

Survival rates for lung cancer vary depending on the stage of the disease at diagnosis. The five-year overall survival rate for lung cancer is 20.5%, while for stage I lung cancer, the survival rate is 59%, and for stage II, it is 31.7% (NCI, 2019). Early detection and treatment are crucial for improving survival outcomes, as the survival rates for advanced stages of lung cancer are significantly lower. Lung cancer is commonly diagnosed among individuals aged 65-74, with a median age of 71 at diagnosis. In 2020, it was estimated that 12.7% of all cancer cases in the United States would be attributed to lung cancer (NCI, 2019).

There are three common tests used for the diagnosis in lung cancer: (i) Low-dose spiral CT scan (LDCT scan), which utilizes a low-dose radiation X-ray machine to obtain detailed images of the internal body areas; (ii) Chest X-ray, which provides images of the organs and bones within the chest; and (iii)

Sputum cytology, a procedure where a sample of coughed-up mucus (sputum) is examined under a microscope to check for the presence of cancer cells (NCI, 2019). The choice of which test or tests to use depends on various factors such as the individual's symptoms, medical history, and the suspected nature of the lung abnormality.

## 1.2 Probability Modeling of Cancer Screening

The disease progression model is a commonly used framework in cancer screening probability modeling, providing a basis for estimating the relevant probabilities and their parameters (Zelen and Feinleib, 1969). In this study, the disease progression model is adopted as the foundation for the cancer screening probability modeling. This model assumes that the development of cancer follows a sequence of three distinct states, denoted as $S_0 \rightarrow S_p \rightarrow S_c$. These states represent different stages of the disease progression.

The initial state $S_0$, refers to the disease-free state. In this state, an individual does not have the disease or has a disease at an early stage that cannot be detected by the screening examination. The next state, $S_p$, represents the preclinical disease state. In this stage, an asymptomatic individual unknowingly possesses the disease that can be detected through a screening examination. This state is crucial in cancer screening as the goal of screening is to detect the tumor in the preclinical state. A graphical illustration of the disease progression model is presented in figure 1.2.1.

Figure 1.2.1: A graphical representation of the disease progression model

Several definitions and notations related to probability modeling are introduced in this chapter, which are utilized throughout the remainder of the dissertation. The focus is on lung cancer screening, specifically for a cohort of initially asymptomatic individuals with no prior history of lung cancer.

One important parameter in this study is sensitivity, which represents the probability of correctly identifying a disease during screening when an individual is in the preclinical state ($S_p$). The sensitivity is estimated using an epidemiological approach, which assumes that sensitivity is not dependent on the age of participants (Wang et al., 2017). To estimate sensitivity, the total number of cases detected through screening is divided by the sum of screen-detected cases and interval cases.

Another key parameter to be estimated is the sojourn time, which refers to the duration that a person remains in the preclinical state ($S_p$). This time period represents the interval during which a person is asymptomatic, but the cancer is detectable through screening. A longer sojourn time implies that the disease can be detected more easily through screening. The sojourn time ($T_p$) is calculated as the difference between the age at which a person enters the preclinical state and the age at which clinical symptoms appear ($T_p = t_2 - t_1$), as illustrated in Figure 1.2.1. The mean sojourn time (MST) is the average time spent in the preclinical screen-detectable phase.

The transition probability is another key parameter estimated in this

study. It represents the probability density function (PDF) that describes the time duration in the disease-free state $(S_0)$. Additionally, it provides essential information about the age at which individuals transition from the disease-free state to the preclinical state $(S_p)$. The transition density is typically estimated using common parametric models or assumptions of constant intervals.

In general, it is assumed that the sojourn time and transition time are independent (Wu et al., 2005). However, since the inception of either $S_p$ or $S_c$ is difficult to observe directly in a screening program, proper modeling is required for accurate estimation of the sojourn time and transition density. The objective of this study is to provide accurate statistical inference for the distribution of the sojourn time and the transition probability from the disease-free to the preclinical state for heavy smokers, utilizing the chest X-ray data from the National Lung Screening Trial (NLST).

The existing statistical methods used to estimate these key parameters in lung cancer screening, as well as the methods for estimating the optimal screening age/time, the lead time distribution and the probability of overdiagnosis are briefly reviewed in this study for various scenarios. The focus is on an asymptomatic heavy smoker of age $t_0$, with no history of lung cancer, who undergoes a series of ordered screening exams at regular intervals.

Initially, we consider an asymptomatic heavy smoker of age $t_0$, with no lung cancer history. We suppose that the person will undergo $K$ ordered screening exams at ages $t_0 < t_1 < ... < t_{K-1}$, where $t_i = t_0 + i$ is for annual screening exams. We define the $i$th screening interval as the interval between the $i$th and the $(i + 1)$th screening exams $(t_{i-1}, t_i)$ where $i = 1, 2, ..., K - 1$. The sensitivity of the screening exam is denoted as $\beta$, and the function $w(t)$ represents the duration in the disease-free state $(S_0)$, it is a sub-PDF since an individual may remain in the disease-free state throughout their lifetime.

Finally, the probability density function of the sojourn time in $S_p$ is denoted as $q(x)$, with the survival function $Q(z) = \int_z^\infty q(x)dx$.

Table 1.2.1: Sample Cancer Screening Data

| age $(t_0)$ | $n_1$ | $s_1$ | $r_1$ | $n_2$ | $s_2$ | $r_2$ | ... | $^a n_k$ | $s_k$ | $r_k$ |
|---|---|---|---|---|---|---|---|---|---|---|
| ... | ... | | ... | | ... | | | | | ... |
| 60 | 1188 | 4 | 1 | 1123 | 3 | 2 | | 1091 | 1 | 2 |
| ... | | | | | | | | | | |
| 65 | 752 | 6 | 1 | 704 | 3 | 3 | | 686 | 3 | 3 |
| ... | ... | | ... | | | | | | | |

$^a$ The total number of screening $K > 0$ is an integer. $K = 3$ in NLST study

In the context of cancer screening, the screening data typically consists of three parts for each screening cycle. Let's define the following variables: (i) $n_{i,t_0}$ represents the total number of individuals in the cohort who undergo the $i$th screening exam. This variable provides the population size at a particular screening point; (ii) $s_{i,t_0}$ denotes the number of cases detected at the $i$th screening exam. These are the individuals who are identified and confirmed as having cancer during the screening process; (iii) $r_{i,t_0}$ refers to the number of cases diagnosed in the clinical state $(S_c)$ within the interval between the $(i-1)$th and $i$th screening exams, these cases are known as interval cases and are diagnosed with cancer between two consecutive screening intervals.

Table 1.2.1 illustrates the data format for a screening program with $K$ scheduled exams. The table showcases the stratification of the data based on the initial age at which the screenings are conducted. The three parts, $n_i, s_i, r_i$ represent the variables discussed above, corresponding to the respective screening exams.

The probability of overdiagnosis is a significant concern in lung cancer screening and the focus of the second part of this dissertation. It occurs when

cancer is detected through screening, but the individual would not exhibit any clinical symptoms. Essentially, overdiagnosis refers to the detection of cancers that would not have caused any medical or clinical problems in the absence of screening. This can lead to unnecessary treatments or interventions for individuals who would not have experienced any adverse health outcomes.

To better understand the factors influencing the risk of overdiagnosis, it is essential to estimate various events prior to the diagnosis status and the appearance of symptoms before death. These events include the concept of symptom-free-life, no-early-detection, true-early-detection, and others. Table 1.2.2 provides a summary of these events, which are relevant for assessing the risk of overdiagnosis.

Table 1.2.2: Definition of long-term outcomes/events in screening

| Diagnosis status | Ultimate lifetime disease status | |
|---|---|---|
| | No symptom before death | symptoms before death |
| not screen detected | symptom free life | no early detection |
| screen detected | overdiagnosis | true early detection |

People who take part in cancer screening are divided into four mutually exclusive groups: true-early-detection, no-early-detection, overdiagnosis, and symptom-free-life. In this dissertation (chapter 3) the probability for each case derived in Wu et al. (2014) is used. These probabilities change with a person's age at study entry, screening frequency, screening sensitivity, and other parameters. It is also allowed that the human lifetime is subject to a competing risk of death from other causes. The model can provide policymakers with important information regarding the distribution of individuals participating in a screening program which eventually fall into one of the four groups:

- *Group 1*: Symptom-free-life (SympF)- The participant went through the screening exams, but lung cancer was never got detected, and ultimately

he/she died of other causes.

- *Group 2*: No-early-detection (NoED)- The participant went through the screening exams, and the disease reveals clinically but was not detected by scheduled screening exams.

- *Group 3*: True-early-detection (TrueED)- The participant was diagnosed with lung cancer at a scheduled screening exam, and his/her clinical symptoms would have appeared before his/her death.

- *Group 4*: Overdiagnosis (OverD)- The participant of the Group 4 was diagnosed with lung cancer at a scheduled screening exam, but his/her clinical symptoms would not appear before his/her death.

To estimate the probability of overdiagnosis, we introduce the time variable $t$, which represents an individual's age at the time of screening. Additionally, we consider the variable $T$, which represents a person's lifetime and is modeled as a continuous random variable with a probability density function $f_T(t)$. We define an event denoted as $A$, which indicates that participants in the screening exams are asymptomatic for cancer before and at the time $t_0$. In other words, event $A$ signifies that individuals do not exhibit any clinical symptoms related to cancer up until the time of their initial screening. This assumption allows us to consider the probability of overdiagnosis in the context of individuals who are initially free of cancer symptoms.

A = { Participant is asymptomatic of lung cancer before and at $t_0$ }

The conditional probability of $A$ derived in Wu et al. (2014) considers that a participant is asymptomatic at $t_0$ (event A), no lung cancer was found before age $t_0$, given that one's lifetime $T$ exceeds $t_0$, is the sum of the two

probabilities: (i) participant remains in the disease-free state through age $t_0$, the probability of which is $1 - \int_0^{t_0} w(x)dx$, and (ii) participant enter preclinical state, $S_p$ before $t_0$ but remains in $S_p$ for a long time that no symptoms appear before $t_0$, the probability of which is $\int_0^{t_0} w(x)Q(t_0 - x)dx$.

$$P(A|T > t_0) = 1 - \int_0^{t_0} w(x)dx + \int_0^{t_0} w(x)Q(t_0 - x)dx \qquad (1.1)$$

The model where a person's lifetime is a random variable, $T \sim f_T(t)$ is applied in this study to estimate the probability of each case considering multiple screening exams.

The method employed in the third part of this study focuses on estimating the optimal screening time and evaluating the lead time, and overdiagnosis. The approach is based on the probability model developed by Wu (2022) and involves the following steps:

- *Optimal screening time/age*: The first step is to find the optimal age for the first screening exam, denoted as $t_0$, for an asymptomatic individual of current age $a_0$. By constraining incidence probability to a small value, such as 10% or 20%, a unique solution for the first screening time can be obtained.

- *Lead time distribution*: The lead time distribution is estimated for individuals who would be diagnosed with cancer at their first screening exam.

- *Probability of overdiagnosis*: The probability of overdiagnosis is calculated given that a person is diagnosed at the first screening exam and their human lifetime exceeds the first screening age.

This methodology allows for optimal scheduling on an individual basis

and provides accurate estimation of lead time, overdiagnosis, and true-early-detection probabilities at the future screening time if one would be diagnosed with cancer.

# CHAPTER 2

# INFERENCE OF SOJOURN TIME AND TRANSITION DENSITY USING THE NLST X-RAY SCREENING DATA IN LUNG CANCER

## 2.1 Introduction

The key parameters in the probability modeling of cancer screening include the sojourn time distribution, the transition probability, and screening sensitivity. These parameters play a crucial role in understanding and evaluating the effectiveness of cancer screening programs. Additionally, other important measures such as the probability of overdiagnosis, optimal screening time, and lead time distribution can be expressed as functions of these key parameters.

Accurate estimation of these key parameters is essential for ensuring the reliability and validity of cancer screening models. In this chapter, the main objective is to provide precise statistical inference for the distribution of sojourn time and the transition probability from the disease-free state to the preclinical state specifically for heavy smokers. This estimation is carried out using the chest X-ray data from the National Lung Screening Trial (NLST). An existing conditional likelihood function Wu et al. (2005) will be used to achieve this goal of evaluating lung cancer screening.

The subsequent sections of this chapter are structured as follows: Section 2.2 presents an introduction to the NLST data, Section 2.3 focuses on estimating

11

key parameters for asymptomatic participants of current age $t_0$, Section 2.4 demonstrates the application of an existing method through simulation utilizing a Bayesian and MCMC approach, and Section 2.5 engages in a discussion regarding the results obtained from the simulation.

## 2.2   The National Lung Screening Trial

The National Lung Screening Trial (NLST) is a randomized clinical trial. It was launched by the National Cancer Institute in 2002 (Jang et al., 2013a). NLST has screened a high-risk population with either low-dose helical (spiral) computed tomography (CT) or standard chest X-ray (X-ray). In this study, the standard chest X-ray (a single image of the whole chest) data was used, which was divided into two groups: males (15,396) and females (10,634) of heavy smokers. Asymptomatic participants aged 55 to 74 from 33 centers across the US between August 2002 and April 2004 were initially screened from each group. Three annual screening exams were provided to each participant from each group. The data were organized in such a way for accurate estimation: for each age $t_0$ at study entry, and each screening, the total number of people being screened $n_i$, the number of confirmed cancer cases $s_i$ and the number of interval cases $r_i$, before the next exam. Participants that dropped in the middle of the program are also included. Table 1.2.1 shows the data format that have been used in the NLST study. Participants with different ages, gender, and smoking status are considered significant risk factors in this project. If any of the tests were positive, the screen was considered positive, and a definitive workup exam, such as a biopsy, was done.

## 2.3 Method

Let $t$ represents the age of participants in the screening. $\beta(t)$ represents the sensitivity of the screening. We define $w(t)dt$ as the transition probability from $S_0$ to $S_p$ in the time interval $(t, t+dt)$. Let $q(x)$ be the probability density function (pdf) of the sojourn time in $S_p$, where $Q(z) = \int_x^\infty q(x)dx$ is the survival function of the sojourn time in $S_p$.

Initially, an asymptomatic heavy smoker of age $t_0$ is considered, who has no history of lung cancer. Suppose that the person will undergo $K$ screening exams at ages $t_0 < t_1 < ... < t_{K-1}$, where $t_i = t_0 + i$ for annual screening exams in the NLST study. We define the $i$th screening interval as the interval between the $i$th and the $(i+1)$th screening exams $(t_{i-1}, t_i)$ where $i = 1, 2, ..., K-1$. Let $t_{-1} = 0$. For each screening exam, let $n_{i,t_0}$ be the total number of individuals in this cohort examined at the $i$th screening, $s_{i,t_0}$ is the number of cases detected at the $i$th screening exam, and $r_{i,t_0}$ is the number of cases diagnosed in the clinical state $S_c$ within the interval $(t_{i-1}, t_i)$, which is the interval cases. For the NLST chest X-ray data, since the age of participants enrolled was between 55 to 74 at the study entry, the likelihood function for all groups is:

$$L = \prod_{t_0=55}^{74} \prod_{k=1}^{3} D_{k,t_0}^{s_k,t_0} I_{k,t_0}^{r_k,t_0} (1 - D_{k,t_0} - I_{k,t_0})^{n_{k,t_0} - s_{k,t_0} - r_{k,t_0}} \qquad (2.2)$$

where $D_{k,t_0}$ is the probability that an individual will be diagnosed at the $k$th scheduled exam given that he/she is in $S_p$, and $I_{k,t_0}$ is the probability of being incident in the $k$th screening interval. These two probabilities were originally derived in Wu et al. (2005) as follows:

$$D_{1,t_0} = \beta(t_0) \int_0^{t_0} w(x) Q(t_0 - x) dx \qquad (2.3)$$

13

$$D_{k,t_0} = \beta(t_{k-1}) \left\{ \sum_{i=0}^{k-2} [1 - \beta(t_i)]...[1 - \beta(t_{k-2})] \int_{t_{i-1}}^{t_i} w(x)Q(t_{k-1} - x)dx \right.$$

$$\left. + \int_{t_{k-2}}^{t_{k-1}} w(x)Q(t_{k-1} - x)dx \right\}, \qquad \text{for all k} = 2, ..., K \tag{2.4}$$

$$I_{k,t_0}(t) = \sum_{i=0}^{k-1} [1 - \beta(t_i)]...[1 - \beta(t_{k-1})] \int_{t_{i-1}}^{t_i} w(x)[Q(t_{k-1} - x) - Q(t_k - x)]dx$$

$$+ \int_{t_{k-1}}^{t_k} w(x)[1 - Q(t_k - x)]dx, \qquad \text{for all k} = 1, ..., K$$

$$\tag{2.5}$$

The sensitivity was estimated by the epidemiologic method using the NLST data:

$$\beta(t) = \beta_0 = \sum_{t_0=55}^{74} \sum_{k=1}^{K} s_{k,t_0} / \sum_{t_0=55}^{74} \sum_{k=1}^{K} s_{k,t_0} + \sum_{t_0=55}^{74} \sum_{k=1}^{K} r_{k,t_0} \tag{2.6}$$

It was obtained by using the total number of screen-detected cases divided by the sum of screen-detected cases and interval cases (Walter and Day, 1983). The transition density follows a log-Normal PDF multiplied by 30%:

$$w(t|\mu, \sigma^2) = \frac{0.3}{\sqrt{2\pi}\sigma t} exp\{-(logt - \mu)^2/(2 * \sigma^2)\} \tag{2.7}$$

The density of sojourn time follows Weibull distribution.

$$q(x|\alpha, \lambda) = \alpha\lambda x^{\alpha-1} exp(-\lambda x^\alpha) \tag{2.8}$$

So, the survival function of the sojourn time has the following form:

$$Q(x) = exp(-\lambda x^\alpha)$$

where $x$ is the sojourn time in the state of $S_p$; $\alpha$ and $\lambda$ are positive

parameters to be estimated.

## 2.4 Application

In this study, the likelihood function (2.2) was employed to calculate the Bayesian estimate of the four unknown parameters $\theta = (\mu, \sigma^2, \alpha, \lambda)$ using the NLST chest X-ray data. To estimate the sensitivity of the screening study, we adopted the epidemiological method, which assumes that sensitivity does not vary with the age of the participants. The estimated values for sensitivity were found to be 0.61 for males and 0.62 for females. This indicates that the sensitivity is slightly higher in females compared to males.

Table 2.4.1: Bayesian Posterior Estimates

| Parameters | Male | | | Female | | |
|---|---|---|---|---|---|---|
| | Mean | Median | SE | Mean | Median | SE |
| $\mu$ | 4.3158 | 4.3156 | 0.0114 | 4.3178 | 4.3176 | 0.0665 |
| $\sigma^2$ | 0.0220 | 0.0218 | 0.0032 | 0.0376 | 0.0374 | 0.0063 |
| $\alpha$ | 2.3274 | 2.1165 | 0.7130 | 2.7634 | 2.7797 | 0.7200 |
| $\lambda$ | 0.4210 | 0.4326 | 0.0587 | 0.4063 | 0.4164 | 0.0665 |
| $MST$ | 1.2849 | 1.2389 | 0.5862 | 1.2330 | 1.2133 | 0.4823 |

In this study, Markov Chain Monte Carlo (MCMC) methodology was employed to draw posterior samples from the target distribution using non-informative Uniform priors. The prior distributions for the parameters were specified as follows: $\mu$ was assigned a Uniform prior (0.1, 5), $\sigma^2$ had a Uniform prior (0.01, 0.99), $\alpha$ was assigned a Uniform prior (0.1, 5), and $\lambda$ had a Uniform prior (0.1, 2). The parameter ranges chosen were based on previous studies (Wu et al., 2011; Liu et al., 2015), which identified the ranges as $4 < \mu < 4.5$, $0.01 < \sigma^2 < 0.05$, $1.5 < \alpha < 4$, and $0.01 < \lambda < 0.5$.

Three simulations were carried out that were overdispersed with respect

to the target distribution. Each simulation was run for 200,000 iterations, with 50,000 burn-in steps. After the burn-in steps, the posteriors were sampled every 300 steps, resulting in 500 posterior samples for the parameter vector $\theta$. The 500 posterior samples from each of the three chains were pooled for the analysis, providing 1,500 posterior samples for $\theta$. The MCMC trace and the posterior density of $\theta$ are plotted using the final 1,500 posterior samples for $\theta$ of two groups: male and female. Figure 2.4.1 and figure 2.4.2 show the MCMC trace plot of the male and female groups, respectively. Figures 2.4.3 and 2.4.4 show the density plots for the two groups, respectively. The posterior estimates of parameters $\theta$ are listed in table 2.4.1. Mean sojourn time (MST) of male is found slightly higher than female. The larger posterior standard errors are indeed expected in Bayesian inference due to the incorporation of prior uncertainty. Furthermore, the relatively low incidence of lung cancer within each age group can also contribute to these larger standard errors.

Figure 2.4.1: The MCMC trace plots of parameters $\theta = (\mu, \sigma^2, \alpha, \lambda)$ using chest X-ray of male group

Figure 2.4.2: The MCMC trace plots of parameters $\theta = (\mu, \sigma^2, \alpha, \lambda)$ using chest X-ray of female group

18

Figure 2.4.3: The posterior density plots of the parameters $\theta = (\mu, \sigma^2, \alpha, \lambda)$ using chest X-ray of male group

Figure 2.4.4: The posterior density plots of the parameters $\theta = (\mu, \sigma^2, \alpha, \lambda)$ using chest X-ray of female group

Figures 2.4.3 and 2.4.4 present the density plots for each parameter, $\theta = (\mu, \sigma^2, \alpha, \lambda)$, with the Bayesian mean and 95% credible interval indicated for both the male and female groups. These plots provide a visual representation of the posterior distribution of the parameters and allow for a better understanding of their estimated values.

Figures 2.4.5 and 2.4.6 display the posterior quantities of the transition probability and sojourn time for each group. The age-dependent transition probability follows a log-normal distribution that is multiplied by 30%. Notably, the transition probability exhibits a unimodal pattern, reaching its maximum around age 73 for males and 72 for females. This finding suggests that the likelihood of transitioning from the disease-free state to the preclinical state is

highest during these ages.

Regarding the sojourn time, the posterior mean sojourn time (MST) is estimated to be 1.28 years for males and 1.23 years for females. The posterior median sojourn time is slightly lower, at 1.23 years for males and 1.21 years for females. The standard error for the sojourn time is 0.58 for males and 0.48 for females, indicating the uncertainty associated with these estimates. Consequently, the Bayesian analysis suggests that the MST for males appears to be longer than that for females, implying that males may have a slightly longer period in the preclinical state.



Figure 2.4.5: Posterior quantities (2.5%, 50%, 97.5%) of transition probabilities



Figure 2.4.6: Posterior quantities (2.5%, 50%, 97.5%) of sojourn time probabilities

## 2.5  Discussion

This research focused on Bayesian estimation due to the challenges associated with Maximum Likelihood Estimation (MLE) in the context of limited screenings and small sample sizes. As highlighted by Wu and Kim (2020) in their review article, MLE tends to provide less accurate estimates when the number of screenings is less than four and the sample size is small for each age group. Even with a large sample size, although the average of MLE may be closer to the true value, the standard error tends to be large. In the case of the NLST chest X-ray data, the number of screenings was limited to three, making it difficult to obtain reliable MLE estimates.

To estimate the sojourn time and the transition density, the likelihood method (2.2) was applied separately to the NLST chest X-ray data for male and female heavy smokers. Markov Chain Monte Carlo (MCMC) was employed in a Bayesian framework to draw posterior samples and obtain accurate estimations. The primary objective was to provide policymakers with reliable estimates of the sojourn time and the age at which individuals transition into the preclinical state. This information can be valuable for making informed decisions regarding the appropriate age to initiate lung cancer screening exams for heavy smokers and determining the optimal frequency for future screenings.

In this study, the epidemiological method was utilized to estimate the sensitivity, which represents the probability of obtaining a positive result on a screening exam given that an individual is in the preclinical state. This method is consistent with the approach used in Wang et al. (2017). The study conducted by Wang et al. (2017) using the PLCO X-ray data reported a sensitivity of 0.65 for male smokers and 0.68 for female smokers. In our analysis of the NLST chest X-ray data, we found the sensitivity to be 0.61 for males

and 0.62 for female heavy smokers, which closely aligns with the results from the PLCO X-ray data.

Furthermore, previous NCI trials have demonstrated an average sensitivity of 0.69 for chest X-ray screenings, with a range of 0.54 to 0.84 (Gavelli and Giampalma, 2000). In the Mayo Lung Project, Wu et al. (2011) estimated a 95% highest posterior density (HPD) interval for sensitivity as (0.72, 0.98) with a posterior mean of 0.89. Additionally, Liu et al. (2015) used the NLST low-dose CT group data and reported a sensitivity of approximately 0.95 across all male and female groups, highlighting the significant improvement in lung cancer screening sensitivity with low-dose CT scans.

The transition from the disease-free state to the preclinical state can occur at different ages and vary between studies. In our analysis, we found that the transition into the preclinical state could occur as early as before age 50 and continue until after age 90, as depicted in Figure 2.4.5. In Wang et al. (2017), the transition probability for male smokers was observed to peak around age 72.5, while in Liu et al. (2015), the peak was around age 70 for both genders. Wu et al. (2011) found a single maximum at age 68 for males in their study, and in the study conducted by Chen et al. (2014) at the Memorial Sloan-Kettering Cancer Center, the transition probability peaked around age 70. In our study, the transition probability from the disease-free to the preclinical state reached its maximum at around age 73 for males and 72 for females. This information aids in determining the appropriate age to initiate screening and target individuals who are more likely to transition into the preclinical state.

The mean sojourn time (MST) is an important parameter in cancer screening where the cancer is detectable but no symptoms have appeared yet. In our study, the estimated MST for male smokers in the NLST X-ray study was 1.28 years, while for female smokers, it was 1.23 years. Although

there is a slight difference, it is not statistically significant. Comparing our results to previous studies, both male and female smokers in our study have a shorter MST than reported in other studies. For example, Liu et al. (2015) estimated the MST to be 1.44 years for males and 1.62 years for females in a CT scan study of lung cancer. In contrast, Chen et al. (2014) found an MST of approximately 3.35 years for male smokers. The Mayo Lung Project study (Wu et al., 2011) reported a shorter MST of 2.24 years. A review article by Chien et al. (2008) summarized MST estimates ranging from 1.38 to 3.86 years from various CT scans. In Wang et al. (2017), the MST for males was 1.50 years, and for females, it was 1.74 years.

The 95% highest posterior density (HPD) interval for the sojourn time in our study was (1.102, 1.576) for males and (1.088, 1.512) for females. Additionally, the 90% HPD interval was (1.108, 1.515) for males and (1.100, 1.453) for females. Overall, our findings suggest that lung cancer screening programs have a relatively short time interval to detect lung cancer due to the relatively short sojourn time.

Obtaining accurate and reliable estimates of key parameters is crucial in cancer screening research. These parameters serve as the foundation for various important aspects of cancer screening, including the lead time distribution, optimal screening time, the probability of overdiagnosis, and the probability of early detection. By accurately estimating these key parameters, we can acquire valuable insights into the effectiveness and potential drawbacks of cancer screening programs.

# CHAPTER 3

# INFERENCE OF LONG TERM OUTCOMES AND OVERDIAGNOSIS IN LUNG CANCER SCREENING

## 3.1 Introduction

Cancer screening may detect cancer at an early stage. By the time symptoms appear, cancer may have begun to spread. If any abnormalities in tissues or cancer are found early, patients may have more choices for treatment and usually have a better prognosis, which could increase survival. According to American Lung Association, lung cancer screening may find 80% of lung cancer at an early stage when it is more curable. 70% of lung cancers are found without screening at a later stage when there is little chance for a cure (ALA, 2022). Therefore, early detection may be useful to reduce the mortality of cancer.

Lung cancer screening has a low radiation exposure risk, similar to other screening tests. The amount of radiation is a little bit higher than the amount women are exposed to through a mammogram (Hendrick, 2010). Screening for lung cancer has the potential to reduce mortality, but it may also detect tumors that would not cause clinical symptoms. Thus there is a need to quantify the long-term outcomes of repeated screening, specially overdiagnosis.

Overdiagnosis is essential for understanding early detection as it refers to a screening exam that detects a disease through a scheduled screening exam, but

the clinical symptoms do not appear before death. Overdiagnosis is a concern in lung cancer screening because new imaging technologies can detect tiny lung nodules (Marcus et al., 2006). Although these nodules are considered to be abnormal, their clinical significance remains uncertain. Therefore, overdiagnosis is one of the key issues to consider when identifying the balance of possible benefits and harms due to cancer screening, as it can lead to treatment that is not necessary.

The long-term effects of continued cancer screening can be evaluated by estimating the probability of each group and the risk of overdiagnosis among the screen-detected cases. There is much research on overdiagnosis based on observational studies, however, the results varied greatly due to a lack of modeling. Wu et al. (2014) have developed a probability method for evaluating the long-term effects of cancer screening. They separated all participants into four mutually exclusive groups: symptom-free-life, no-early-detection, true-early-detection, and overdiagnosis. These probabilities are influenced by factors such as age at study entry, screening frequency, screening sensitivity, and other parameters. Table 1.2.2 shows the definition of overdiagnosis and other groups based on diagnosis and disease status.

This chapter begins with a brief review of the probability methods used to estimate overdiagnosis and related probabilities in Section 3.2. A simulation study is then conducted in Section 3.3 to assess the performance of the probability models. Bayesian inference techniques are applied to estimate the conditional probabilities in Section 3.4, utilizing the available data. The results obtained from the probability models are presented in Section 3.5, providing insights into the probabilities of overdiagnosis and other related outcomes. Finally, a discussion of the findings and their implications is provided in Section 3.6.

## 3.2 Method

I will briefly review the probability methods developed in Wu et al. (2014), that will be applied to the NLST chest X-ray data in this chapter for multiple screening exams. Let, an initially asymptomatic individual undergoes $K$ screening exams, occurring at ages $t_0 < t_1 < ... < t_{K-1}$. The conditional probability of a case in any one of the four groups, given that his/her lifetime is $T = t_K (> t_{K-1})$, can be calculated as follows:

For individuals in Group 1 (Symptom-Free), where clinical lung cancer does not occur during their lifetime, the conditional probability is:

$$P(Case1, A|T = t_k) = 1 - \int_0^{t_K} w(x)dx + \int_{t_{K-1}}^{t_K} w(x)Q(t_K - x)dx$$
$$+ \sum_{j=0}^{K-1}(1-\beta_j)...(1-\beta_{K-1}) \int_{t_{j-1}}^{t_j} w(x)Q(t_K - x)dx \tag{3.9}$$

where the conditional probability of a symptom-free case can arise from any one of $(K + 2)$ mutually exclusive events: (a) the person remained in the disease-free state $S_0$ throughout his/her lifetime, the probability of which is $1 - \int_0^{t_K} w(x)dx$, (b) he/she entered the preclinical state $S_p$ when the person was between ages $t_{j-1}$ and $t_j$ ; $j = 0, ..., K - 1$, was not detected by the following $K - 1$ exams, no symptom appeared before his/her death, (c) the person entered $S_p$ after $t_{K-1}$ with no symptoms before his/her death.

A Group 2 case (no-early-detection) where the probability of no-early-detection is calculated from $I_{K,j}$ as the probability of being an interval case in the time interval $(tj - 1, tj)$ in the order of $K$ screening exams:

$$P(Case2, A|T = t_k) = I_{K,1} + I_{K,2} + ... + I_{K,K} \tag{3.10}$$

where $I_{K,j}$, the probability of an interval case in $(t_{j-1}, t_j)$ can be

calculated (Wu et al. (2007)):

$$I_{K,j} = \sum_{i=0}^{j-1}(1-\beta_i)...(1-\beta_{j-1})\int_{t_{i-1}}^{t_i} w(x)[Q(t_{j-1}-x)-Q(t_j-x)]dx$$
$$+ \int_{t_{j-1}}^{t_j} w(x)[1-Q(t_j-x)]dx, \text{ for all j=1,...,K.} \tag{3.11}$$

A Group 3 case, the probability of true-early-detection is:

$$P(Case3, A|T=t_K) = \sum_{j=1}^{K-1}\beta_j\Big\{\sum_{i=0}^{j-1}(1-\beta_i)...(1-\beta_{j-1})\int_{t_{i-1}}^{t_i} w(x)[Q(t_j-x)$$
$$- Q(t_K-x)]dx + \int_{t_{j-1}}^{t_j} w(x)[Q(t_j-x)-Q(t_K-x)]dx\Big\}$$
$$+ \beta_0\int_0^{t_0} w(x)[Q(t_0-x)-Q(t_K-x)]dx$$

$$\tag{3.12}$$

where one of $K$ mutually exclusive events can arises depending on his/her age at diagnosis: if a participant is diagnosed at $t_j$, $j = 0, 1, ..., K-1$, then he/she must have entered the preclinical state $S_p$ before $t_j$ , and missed the previous exams, and his/her sojourn time must have been at least $(t_j - x)$ and at most $(t_K - x)$, where $x$ represent the onset time of the preclincal state.

A Group 4 case, overdiagnosis, also can arise as one of $K$ mutually exclusive events, where a person might be diagnosed at the $j$th exam, but his/her symptoms did not appear before his/her death. Hence, the conditional probability of overdiagnosis is:

$$P(Case4, A|T=t_K) = \sum_{j=1}^{K-1}\beta_j\Big\{\sum_{i=0}^{j-1}(1-\beta_i)...(1-\beta_{j-1})\int_{t_{i-1}}^{t_i} w(x)Q(t_K-x)dx$$
$$+ \int_{t_{j-1}}^{t_j} w(x)Q(t_K-x)dx\Big\} + \beta_0\int_0^{t_0} w(x)Q(t_K-x)dx$$

$$\tag{3.13}$$

And it was verified in Wu et al. (2014) that for any screening number

$K \geq 1$,

$$\sum_{i=1}^{4} P(Casei, A|T = t_K) = 1 - \int_0^{t_0} w(x)dx + \int_0^{t_0} w(x)Q(t_0 - x)dx \tag{3.14}$$

$$= P(A|T \geq t_0)$$

For a future screening schedule, such as $t_0 < t_1 < ...$, considering the screening number $K = K(T)$ is a random variable, which changes with the lifetime $T$, the probability of each case when his/her lifetime $T$ is longer than $t_0$ can be estimated from the following:

$$P(\text{Case i}, A|T \geq t_0) = \int_{t_0}^{\infty} P(\text{Case i}, A|K = K(T), T = t)f_T(t|T \geq t_0)dt,$$

for all i=1,2,3,4

$$\tag{3.15}$$

where $f_T(t|T \geq t_0)$ was defined in Wu et al. (2012) as follows:

$$f_T(t|T \geq t_0) = \begin{cases} \frac{f_T(t)}{P(T>t_0)} = \frac{f_T(t)}{1-F_T(t_0)} & t \geq t_0 \\ 0 & \text{otherwise} \end{cases} \tag{3.16}$$

While evaluating the probabilities of each case, it is necessary to verify that for any future screening schedule when the lifetime $T$ is random, the sum of these probabilities is one (3.17).

$$\sum_{i=1}^{4} P(Casei|A, T \geq t_0) = 1 \tag{3.17}$$

A different sojourn time distribution (Weibull distribution) and sensitivity of age (epidemiological approach) than Wu et al. (2014) is used for the mathematical simplicity, where Weibull is considered more flexible as $n$th moments exist, and the sensitivity is independent of age.

## 3.3  Simulation Study

Simulation studies were conducted using the established method described in Section 3.2. The probability of each case was a function of age at the initial screening, the screening interval, the sensitivity, the sojourn time in the preclinical state, the transition probability from the disease-free to the preclinical state, and the human lifetime, $T$. To determine the effects of these factors on the probability of each outcome and to explore how the proportion of true-early-detection and overdiagnosis change among the screen-detected cases due to these factors, the following scenarios were considered for the simulation: age at initial screening $t_0 =$ 55, 60, 65, screening interval, $\Delta = 12, 18, 24$ months, screening sensitivity, $\beta$ is 0.62 (female) and 0.61 (male) from Rahman and Wu (2021).

The transition probability density was chosen to be a log-Normal pdf, with an upper limit of 30% for lung cancer. The parameters $(\mu, \sigma^2)$ is chosen to be $(4.250, 0.015)$, so that the mode of the transition density is about 70 years old. The sojourn time distribution was chosen to be a Weibull PDF, with parameters, $\alpha = 2.50, 2.00, 1.60$ and $\lambda = 0.30, 0.03, 0.02$ such that the MST is 2, 5, and 10 years, respectively.

The number of screens were considered from $K = K(T) = \lceil \frac{(T-t_0)}{\Delta} \rceil$, where $K$ is the largest integer that is less than or equal to $\frac{(T-t_0)}{\Delta}$ is a function of lifetime $T$. The actuarial life table (2016) was chosen from the Social Security Administration (SSA) for the lifetime distribution (NIH, 2020). The period life table was based on mortality, and it provides the probability of death from age 0 to 119 for both males and females. The conditional lifetime distribution $f_T(t|T \geq t_0)$ was estimated using the life table of 2016. The conditional density function of the lifetime $T$ for males and females were plotted in figure 3.3.1 for

three initial ages at screening, $t_0$ =55, 60, 65; irrespective of screening or any specific causes of death.

In tables 3.3.1 and 3.3.2, the first column, mean sojourn time (MST) is given in years. The following four columns are the conditional probabilities (in percentage) of each of the four cases, i.e., $P(Casei|A, T \geq t_0)$, $i = 1, 2, 3, 4$ corresponding to the probability of symptom-free-life, no-early-detection, true-early-detection, and overdiagnosis. Therefore, the summation of these four probabilities in each row is close to 1. The last two columns are the conditional probability of true-early-detection and overdiagnosis, given that it is a screen-diagnosed case; which is calculated by

$$\frac{P(Casei|A, T \geq t_0)}{P(Case3|A, T \geq t_0) + P(Case4|A, T \geq t_0)}, i = 3, 4.$$

The probabilities are reported as percentages in these tables.

In this simulation, mean sojourn time plays the most important role in the case of overdiagnosis. In the last column of table 3.3.1 and 3.3.2, the proportion of overdiagnosis is as high as 29.03% among the screen-diagnosed cases if the mean sojourn time was 10 years long for male and 29.34% if female. It is around 12.68% for male and 11.40% for female if the mean sojourn time is 5 years long, and it is only about 2% both for male and female when the mean sojourn time changes to 2 years when the age at screening is 65 years with a screening interval of 24 months.

The screening interval also plays a role in these probabilities: when the screening interval is longer, the probability of no-early-detection is larger, the probability of true-early-detection is smaller, and the probability of overdiagnosis is slightly smaller. The case of symptom-free-life is pretty stable in all the simulations, it is about 72-77% for the whole population.

Table 3.3.1: Simulation: when the sojourn time is Weibull (Male)

| MST | $P^a$(SympF) | $P^a$(NoED) | $P^a$(TrueED) | $P^a$(OverD) | $P^b$(TrueED\|ScrD) | $P^b$(OverD\|ScrD) |
|---|---|---|---|---|---|---|
| | | | **Screening interval** $\Delta = 12month$, $t_0 = 55$, $\beta = 0.61$ | | | |
| 2 | 76.71 | 9.32 | 13.04 | 0.33 | 97.52 | 2.47 |
| 5 | 76.79 | 7.43 | 12.81 | 2.59 | 83.17 | 16.82 |
| 10 | 76.77 | 5.82 | 10.48 | 6.22 | 62.77 | 37.22 |
| | | | **Screening interval** $\Delta = 18month$, $t_0 = 55$, $\beta = 0.61$ | | | |
| 2 | 76.70 | 10.75 | 12.11 | 0.34 | 97.25 | 2.75 |
| 5 | 76.93 | 7.58 | 12.14 | 2.50 | 82.90 | 17.09 |
| 10 | 76.96 | 5.91 | 10.17 | 6.12 | 62.46 | 37.54 |
| | | | **Screening interval** $\Delta = 24month$, $t_0 = 55$, $\beta = 0.61$ | | | |
| 2 | 75.55 | 12.38 | 10.93 | 0.23 | 97.95 | 2.05 |
| 5 | 76.07 | 7.79 | 13.33 | 2.16 | 86.07 | 13.93 |
| 10 | 76.19 | 6.01 | 11.82 | 5.73 | 67.35 | 32.65 |
| | | | **Screening interval** $\Delta = 12month$, $t_0 = 60$, $\beta = 0.61$ | | | |
| 2 | 76.36 | 9.19 | 13.87 | 0.38 | 97.35 | 2.65 |
| 5 | 75.69 | 7.58 | 13.09 | 2.71 | 82.83 | 17.17 |
| 10 | 75.38 | 6.02 | 11.89 | 6.49 | 64.69 | 35.31 |
| | | | **Screening interval** $\Delta = 18month$, $t_0 = 60$, $\beta = 0.61$ | | | |
| 2 | 76.36 | 10.57 | 12.01 | 0.38 | 96.96 | 3.04 |
| 5 | 75.87 | 7.73 | 13.41 | 2.60 | 83.75 | 16.25 |
| 10 | 75.61 | 6.10 | 11.58 | 6.37 | 64.52 | 35.48 |
| | | | **Screening interval** $\Delta = 24month$, $t_0 = 60$, $\beta = 0.61$ | | | |
| 2 | 74.88 | 12.13 | 11.90 | 0.27 | 97.77 | 2.23 |
| 5 | 74.68 | 7.93 | 14.59 | 2.26 | 86.61 | 13.39 |
| 10 | 74.52 | 6.20 | 13.22 | 5.98 | 68.87 | 31.13 |
| | | | **Screening interval** $\Delta = 12month$, $t_0 = 65$, $\beta = 0.61$ | | | |
| 2 | 75.83 | 8.22 | 15.42 | 0.41 | 97.38 | 2.62 |
| 5 | 73.74 | 7.29 | 15.39 | 2.77 | 84.77 | 15.23 |
| 10 | 72.57 | 6.01 | 13.82 | 6.71 | 67.33 | 32.67 |
| | | | **Screening interval** $\Delta = 18month$, $t_0 = 65$, $\beta = 0.61$ | | | |
| 2 | 75.82 | 9.44 | 13.98 | 0.43 | 96.99 | 3.01 |
| 5 | 73.89 | 7.44 | 15.78 | 2.72 | 85.29 | 14.70 |
| 10 | 72.76 | 6.08 | 13.54 | 6.67 | 66.99 | 33.01 |
| | | | **Screening interval** $\Delta = 24month$, $t_0 = 65$, $\beta = 0.61$ | | | |
| 2 | 74.73 | 10.79 | 14.13 | 0.32 | 97.84 | 2.16 |
| 5 | 73.09 | 7.60 | 16.02 | 2.33 | 87.32 | 12.68 |
| 10 | 72.07 | 6.16 | 15.19 | 6.21 | 70.97 | 29.03 |

[a] The probability of each outcomes, i.e., $P(Casei|A, T \geq t_0)$
[b] The conditional probability of True-Early-Detection and of Over-Diagnosis given that it is a screen-diagnosed

Table 3.3.2: Simulation: when the sojourn time is Weibull (Female)

| MST | $P^a$(SympF) | $P^a$(NoED) | $P^a$(TrueED) | $P^a$(OverD) | $P^b$(TrueED\|ScrD) | $P^b$(OverD\|ScrD) |
|---|---|---|---|---|---|---|
| | | **Screening interval** $\Delta = 12month$, $t_0 = 55$, $\beta = 0.62$ | | | | |
| 2 | 75.32 | 10.06 | 13.72 | 0.26 | 98.16 | 1.84 |
| 5 | 75.34 | 8.19 | 14.11 | 2.28 | 86.09 | 13.90 |
| 10 | 75.31 | 6.62 | 11.78 | 5.92 | 66.54 | 33.46 |
| | | **Screening interval** $\Delta = 18month$, $t_0 = 55$, $\beta = 0.62$ | | | | |
| 2 | 75.31 | 11.66 | 12.45 | 0.27 | 97.90 | 2.09 |
| 5 | 75.47 | 8.37 | 13.37 | 2.20 | 85.85 | 14.15 |
| 10 | 75.48 | 6.72 | 11.44 | 5.84 | 66.21 | 33.79 |
| | | **Screening interval** $\Delta = 24month$, $t_0 = 55$, $\beta = 0.62$ | | | | |
| 2 | 74.64 | 13.49 | 11.02 | 0.17 | 98.47 | 1.53 |
| 5 | 75.04 | 8.61 | 13.47 | 1.91 | 87.59 | 12.41 |
| 10 | 75.14 | 6.84 | 12.05 | 5.51 | 68.62 | 31.38 |
| | | **Screening interval** $\Delta = 12month$, $t_0 = 60$, $\beta = 0.62$ | | | | |
| 2 | 75.59 | 9.77 | 14.30 | 0.28 | 98.06 | 1.94 |
| 5 | 74.89 | 8.22 | 14.12 | 2.35 | 85.76 | 14.24 |
| 10 | 74.57 | 6.72 | 11.98 | 6.09 | 66.31 | 33.69 |
| | | **Screening interval** $\Delta = 18month$, $t_0 = 60$, $\beta = 0.62$ | | | | |
| 2 | 75.59 | 11.31 | 12.17 | 0.28 | 97.73 | 2.27 |
| 5 | 75.04 | 8.39 | 13.39 | 2.25 | 85.59 | 14.40 |
| 10 | 74.75 | 6.81 | 11.64 | 5.98 | 66.06 | 33.94 |
| | | **Screening interval** $\Delta = 24month$, $t_0 = 60$, $\beta = 0.62$ | | | | |
| 2 | 74.71 | 13.05 | 11.83 | 0.19 | 98.39 | 1.61 |
| 5 | 74.39 | 8.61 | 14.51 | 1.97 | 88.07 | 11.93 |
| 10 | 74.19 | 6.92 | 12.25 | 5.66 | 68.41 | 31.59 |
| | | **Screening interval** $\Delta = 12month$, $t_0 = 65$, $\beta = 0.62$ | | | | |
| 2 | 76.29 | 8.67 | 14.60 | 0.31 | 97.91 | 2.09 |
| 5 | 74.13 | 7.80 | 15.12 | 2.38 | 86.40 | 13.59 |
| 10 | 72.95 | 6.59 | 13.68 | 6.22 | 68.74 | 31.26 |
| | | **Screening interval** $\Delta = 18month$, $t_0 = 65$, $\beta = 0.62$ | | | | |
| 2 | 76.28 | 9.98 | 12.90 | 0.33 | 97.53 | 2.47 |
| 5 | 74.25 | 7.94 | 15.44 | 2.33 | 86.87 | 13.13 |
| 10 | 73.11 | 6.67 | 13.35 | 6.18 | 68.37 | 31.63 |
| | | **Screening interval** $\Delta = 24month$, $t_0 = 65$, $\beta = 0.62$ | | | | |
| 2 | 75.56 | 11.48 | 11.87 | 0.22 | 98.14 | 1.86 |
| 5 | 73.78 | 8.13 | 15.63 | 2.01 | 88.59 | 11.40 |
| 10 | 72.73 | 6.76 | 13.98 | 5.80 | 70.66 | 29.34 |

[a] The probability of each outcomes, i.e., $P(Casei|A, T \geq t_0)$
[b] The conditional probability of True-Early-Detection and of Over-Diagnosis given that it is a screen-diagnosed

Figure 3.3.1: The conditional PDF of the lifetime of males and females derived from the life table when $t_0 = 55, 60, 65$

## 3.4 Bayesian inference using the NLST data

We applied the existing methods reviewed in section 3.2 to the NLST chest X-ray data. The probability for each of the four cases was the function of the three key parameters, $\beta(t)$, $w(t)$, and $q(x)$. These key parameters were estimated from the NLST chest X-ray data using the models described in section 2.3. For more detail, explore Rahman and Wu (2021).

The posterior predictive probability of each case was estimated using the NLST data from the following:

$$P(Casei|T \geq t_0, A, NLST) = \int P(Casei, \theta|T \geq t_0, A, NLST)d\theta$$
$$= \int P(Casei|T \geq t_0, A, \theta)f(\theta|NLST)d\theta$$
$$\approx \frac{1}{n}\sum_{j=1}^{n} P(Casei|T \geq t_0, A, \theta_j^*).$$

(3.18)

Where $\theta_j^*$ are posterior samples from the MCMC simulation, and the posterior sample size $n = 1500$.

## 3.5 Results

The MCMC posterior samples of 1500 were used in equation 3.18, to conduct Bayesian inference for three hypothetical cohorts of asymptomatic participants. In the first screening exam, the three cohorts had initial screening ages of 55, 60, and 65. For each group, various screening frequencies were examined, and screening intervals of $\Delta = 12, 18, 24$ months, respectively. The number of screens, $K = K(T) = \lceil \frac{(T-t_0)}{\Delta} \rceil$ was considered a function of the lifetime $T$ again. The conditional lifetime distribution (equation 3.16) was estimated using the actuarial life table from the SSA (NIH, 2020) described in the previous section. The conditional probabilities of each of the four cases $P(Casei|A; T \geq t_0; NLST)$ are reported in table 3.5.1 and 3.5.2.

Across all three age groups, the probability of overdiagnosis is observed to be very low. Specifically, for the 12-month screening interval and initial screening ages of 55, 60, and 65, the probabilities of overdiagnosis are approximately 0.33%, 0.38%, and 0.43% for males, and 0.23%, 0.24%, and 0.29% for females, respectively. It is worth noting that these probabilities decrease as the screening interval ($\Delta$) increases. While the probability of overdiagnosis is slightly higher when the initial screening age is 65, there is minimal difference observed for the other age groups. The results presented in Table 3.5.1 and 3.5.2 demonstrate that males exhibit a higher susceptibility to overdiagnosis compared to females.

The probability of true-early-detection for different initial screening ages of 55, 60, and 65 is determined to be 10.72%, 10.32%, and 10.11% for males, and 11.10%, 10.76%, and 9.81% for females, respectively, when annual screening is conducted. It is noteworthy that this probability remains relatively stable as the age at the initial screening exam increases. However, the probability of

35

Table 3.5.1: A projection of lung cancer screening effects using the NLST chest X-ray data (Male)

| $\Delta^a$ | $P^b$(SympF) | P(NoED) | P(TrueED) | P(OverD) |
|---|---|---|---|---|
| **Age at initial screen $t_0 = 55$** | | | | |
| 12 mo. | 80.49 (0.76) | 8.36 (0.40) | 10.72 (0.57) | 0.33 (0.06) |
| 18 mo. | 80.53 (0.76) | 9.31 (0.49) | 9.82 (0.43) | 0.32 (0.08) |
| 24 mo. | 80.62 (0.77) | 11.03 (0.66) | 8.07 (0.40) | 0.23 (0.05) |
| **Age at initial screen $t_0 = 60$** | | | | |
| 12 mo. | 80.84 (0.82) | 8.25 (0.42) | 10.32 (0.59) | 0.38 (0.07) |
| 18 mo. | 80.97 (0.83) | 9.48 (0.52) | 9.12 (0.47) | 0.37 (0.08) |
| 24 mo. | 81.01 (0.83) | 11.31 (0.68) | 7.60 (0.43) | 0.24 (0.06) |
| **Age at initial screen $t_0 = 65$** | | | | |
| 12 mo. | 81.42 (0.82) | 8.02 (0.42) | 10.11 (0.58) | 0.43 (0.07) |
| 18 mo. | 81.53 (0.82) | 9.56 (0.49) | 8.35 (0.47) | 0.41 (0.09) |
| 24 mo. | 81.71 (0.82) | 11.47 (0.64) | 6.43 (0.44) | 0.32 (0.06) |

[a] $\Delta = t_i - t_{i-1}$ is the time interval between screens
[b] The mean probability and it's standard error (in parenthesis) are reported as percentages in the table

true-early-detection decreases with longer screening time intervals. Table 3.5.1 and 3.5.2 provide a comparison, clearly indicating that males have a slightly lower propensity for true-early-detection compared to females.

The probability of no-early-detection is 8.36%, 8.25%, and 8.02% for males, and 7.99%, 7.69%, and 6.91% for females, when participants initiate screening at ages 55, 60, and 65, respectively, using a 12-month screening schedule. This probability increases as the screening interval increases and decreases as the age at initial screening increases.

The probability of symptom-free-life is very high for all age groups, ranging from approximately 80% to 82% for males and around 80% to 83% for females. This probability increases as the screening interval increases. Notably, the difference between the corresponding probabilities is smaller among the age groups of 55 to 65.

Boxplots of the results for the probabilities of each case when $t_0 = 55$, 60, and 65 for both males and females are presented in figures 3.5.1 to 3.5.6.

Table 3.5.2: A projection of lung cancer screening effects using the NLST chest X-ray data (Female)

| $\Delta^a$ | $P^b$(SympF) | P(NoED) | P(TrueED) | P(OverD) |
|---|---|---|---|---|
| | **Age at initial screen $t_0 = 55$** | | | |
| 12 mo. | 80.64 (0.97) | 7.99(0.49) | 11.10 (0.68) | 0.23 (0.06) |
| 18 mo. | 80.68 (0.98) | 9.69 (0.63) | 8.89 (0.56) | 0.22 (0.05) |
| 24 mo. | 80.79 (0.98) | 11.15 (0.79) | 7.91 (0.51) | 0.13 (0.03) |
| | **Age at initial screen $t_0 = 60$** | | | |
| 12 mo. | 81.28 (1.02) | 7.69 (0.51) | 10.76 (0.70) | 0.24 (0.06) |
| 18 mo. | 81.67 (1.01) | 9.16 (0.63) | 8.64 (0.58) | 0.23 (0.05) |
| 24 mo. | 81.97 (1.01) | 10.54 (0.79) | 6.98 (0.52) | 0.17 (0.04) |
| | **Age at initial screen $t_0 = 65$** | | | |
| 12 mo. | 82.69 (0.97) | 6.91 (0.48) | 9.81 (0.67) | 0.29(0.06) |
| 18 mo. | 82.86 (0.96) | 8.88 (0.58) | 7.91 (0.56) | 0.26 (0.05) |
| 24 mo. | 82.97 (0.97) | 9.78 (0.72) | 6.67 (0.50) | 0.20 (0.04) |

[a] $\Delta = t_i - t_{i-1}$ is the time interval between screens
[b] The mean probability and its standard error (in parenthesis) are reported as percentages in the table

In figure 3.5.1, for the age group of 55 in males, the probabilities of symptom-free-life and overdiagnosis either remain stable or exhibit minimal changes with the screening time interval. However, the probability of no-early detection increases monotonically with the screening time interval, while the probability of true-early-detection decreases monotonically with the length of the screening time interval. This pattern is consistent across all screening ages for both males and females.

Figure 3.5.1: The boxplot of the estimated probability of male for each case with $t_0 = 55$

Figure 3.5.2: The boxplot of the estimated probability of male for each case with $t_0 = 60$

Figure 3.5.3: The boxplot of the estimated probability of male for each case with $t_0 = 65$

Figure 3.5.4: The boxplot of the estimated probability of female for each case with $t_0 = 55$

Figure 3.5.5: The boxplot of the estimated probability of female for each case with $t_0 = 60$

Figure 3.5.6: The boxplot of the estimated probability of female for each case with $t_0 = 65$

Table 3.5.3: The estimated probability of male and female given that it is a diagnosed cancer case

| | Male | | | Female | | |
|---|---|---|---|---|---|---|
| $\Delta$ | $P^c(\text{NoED}|D^d)$ | $P(\text{TrueED}|D)$ | $P(\text{OverD}|D)$ | $P(\text{NoED}|D)$ | $P(\text{TrueED}|D)$ | $P(\text{OverD}|D)$ |
| **Age at initial screen $t_0 = 55$** | | | | | | |
| 12 mo. | 42.19 | 56.20 | 1.61 | 40.96 | 57.91 | 1.13 |
| 18 mo. | 52.19 | 46.23 | 1.58 | 51.56 | 47.27 | 1.08 |
| 24 mo. | 60.84 | 37.85 | 1.31 | 61.25 | 38.01 | 0.74 |
| **Age at initial screen $t_0 = 60$** | | | | | | |
| 12 mo. | 42.05 | 56.01 | 1.94 | 40.86 | 57.87 | 1.19 |
| 18 mo. | 51.82 | 46.41 | 1.77 | 51.34 | 47.49 | 1.17 |
| 24 mo. | 60.39 | 38.06 | 1.55 | 60.91 | 38.20 | 0.89 |
| **Age at initial screen $t_0 = 65$** | | | | | | |
| 12 mo. | 41.84 | 55.87 | 2.29 | 40.72 | 57.78 | 1.50 |
| 18 mo. | 51.10 | 46.80 | 2.10 | 50.94 | 47.72 | 1.34 |
| 24 mo. | 59.67 | 38.44 | 1.89 | 60.42 | 38.42 | 1.16 |

[c] The estimated conditional probability was calculated as $p_i^*/(p_2^* + p_3^* + p_4^*)$, $i = 2, 3, 4$, for each of the 1500 posterior samples, then averaged. It is in percentage.

[d] The event $D=\{$ Diagnosed cases: including both interval incident and screen-detected cases $\}$

The conditional probabilities were evaluated for cases 2, 3, and 4, given that the person was diagnosed with cancer. The probabilities of overdiagnosis are 1.61%, 1.94%, and 2.29% for males, and 1.13%, 1.19%, and 1.50% for females, in the 12-month screening group when the starting age is 55, 60, and 65, respectively. The conditional probability of true-early-detection, given a diagnosed case, decreases significantly as the screening interval $\Delta$ increases. It ranges from around 56% to 38% in the 60-year-old male group and from around 58% to 38% in the 60-year-old female group. Conversely, the conditional probability of no-early-detection increases within each age group as the screening interval increases. These results are summarized in Table 3.5.3.

The probabilities and 95% HPD intervals of true-early-detection and overdiagnosis, given it is a screen-detected case, are listed in table 3.5.4. The length of the 95% HPD interval for these two probabilities (percentages) decreases as the screening interval increases. However, these credible interval lengths increase as the initial screening age increases. For males, the percentage of overdiagnosis is 3.87%, 4.26%, and 5.12% for different screening ages with a 24-month screening interval, which is higher than any other age group and screening interval. For females, the percentages are 2.43%, 2.71%, and 3.46%. In summary, the probability of overdiagnosis is much lower than expected, while the probability of true-early-detection is often above 94% and higher.

Table 3.5.4: The estimated probability of male and female for the screen detected cases with 95% credible interval

| | Male | | Female | |
|---|---|---|---|---|
| $\Delta$ | $P^d(\text{TrueED}|ScrD^e)$ | $P(\text{OverD}|ScrD)$ | $P(\text{TrueED}|ScrD)$ | $P(\text{OverD}|ScrD)$ |
| **Age at initial screen $t_0 = 55$** | | | | |
| 12 mo. | 97.06 (95.79,98.32) | 2.94 (1.68,3.91) | 98.25 (97.11,98.96) | 1.75 (1.04,2.89) |
| 18 mo. | 96.69 (95.54,97.66) | 3.31 (2.34,4.46) | 98.09 (97.11,98.69) | 1.90 (1.31,2.89) |
| 24 mo. | 96.13 (94.84,97.31) | 3.87 (2.69,5.16) | 97.57 (96.42,98.26) | 2.43 (1.74,3.58) |
| **Age at initial screen $t_0 = 60$** | | | | |
| 12 mo. | 96.65 (95.30,97.94) | 3.35 (2.06,4.69) | 97.97 (96.72,98.71) | 2.03 (1.29,3.27) |
| 18 mo. | 96.12 (94.92,97.08) | 3.88 (2.92,5.07) | 97.73 (96.59,98.28) | 2.27 (1.71,3.41) |
| 24 mo. | 95.74 (94.39,96.89) | 4.26 (3.10,5.61) | 97.29 (96.03,97.97) | 2.71 (2.02,3.97) |
| **Age at initial screen $t_0 = 65$** | | | | |
| 12 mo. | 96.06 (94.63,97.46) | 3.94 (2.53,5.36) | 97.46 (96.08,98.31) | 2.54 (1.68,3.91) |
| 18 mo. | 95.34 (94.13,96.36) | 4.66 (3.63,5.87) | 97.08 (95.84,97.72) | 2.92 (2.28,4.16) |
| 24 mo. | 94.89 (93.44,96.09) | 5.12 (3.91,6.56) | 96.54 (95.12,97.32) | 3.46 (2.67,4.87) |

[d] The estimated conditional probability was calculated as $p_i^*/(p_3^* + p_4^*)$, $i = 3, 4$, for each of the 1500 posterior samples, then averaged. It is in percentage.

[e] The event $ScrD=\{$ Screen-detected case$\}$

## 3.6   Discussion

In this study, we aimed to assess the long-term effects of lung cancer screening using chest X-rays. Asymptomatic participants in the screening program were eventually separated into four distinct groups: symptom-free-life, no-early-detection, true-early-detection, and overdiagnosis based on their diagnosis status and disease status. Our analyses provide policymakers with valuable estimates of the probability of true-early-detection, overdiagnosis, and other relevant outcomes that arise from a periodic lung cancer screening program. To address uncertainty and calculate variations, we limited our analysis to the Bayesian approach, which allows for the determination of credible intervals (percentages).

On March 9, 2021, the U.S. Preventive Services Task Force released (USPSTF, 2021) new recommendations endorsing annual screening for lung cancer using low-dose computed tomography (LDCT) for individuals aged 50 to 80 years who have a 20 pack-year smoking history and currently smoke or have quit within the past 15 years. Based on data from the NLST chest X-ray, our findings indicate a high probability of symptom-free-life, ranging from 78% to 82% for all participants, i.e., heavy smokers. Additionally, the probability of overdiagnosis is very low, less than 0.43%, among all participants, regardless of gender (as shown in tables 3.5.1 and 3.5.2) for annual screening. The estimated rates of overdiagnosis in the 55-year-old cohort are as follows: for males, the rates are 2.94% if screenings are conducted annually and 3.87% if screenings occur every other year (as shown in Table 9). These estimates come with a 95% Highest Posterior Density (HPD) interval ranging from 2.69% to 5.16%. Comparatively, for females in the 55-year-old cohort, the rates of overdiagnosis among screen-detected cases are 1.75% if screenings take place annually and

2.43% if screenings happen every other year. It's important to note that the overdiagnosis rates are lower for females than for males. Additionally, as individuals age, the probability of overdiagnosis increases for both males and females, even among those cases detected through screening.

The probability of symptom-free-life can be calculated as the complement of the lifetime risk. According to Villeneuve and Mao (1994), the lifetime risk of developing lung cancer is 17.2% for male current smokers and 11.6% for female current smokers. However, the risk is substantially reduced for individuals who have never smoked regularly. For the 55-year-old age group, the estimated probability of symptom-free-life is approximately 80.49% for males and 80.64% for females, which aligns with the accepted lifetime risk. This probability indicates the likelihood of remaining free of symptoms throughout one's life.

In a study by Tammemaegi et al. (2014), it was found that 26% of individuals selected for screening based on USPSTF criteria had risks below the threshold defined by the $PLCO_m2012risk$ lung cancer risk prediction model. This model, derived from the Prostate Lung Colorectal and Ovarian Screening (PLCO) study, assesses the risk of lung cancer based on various factors. For former smokers with a quit time of more than 15 years, 8.5% had risks exceeding the threshold. Notably, the risks of lung cancer were significantly higher in PLCO smokers aged 65-80 years compared to those aged 55-64 years.

In the study by Wu et al. (2016) , which analyzed NLST CT scan data for lung cancer, they found that the probability of overdiagnosis increases with age. Specifically, the probability ranged from 3% to 9% when individuals aged from 60 years to 80 years. This suggests that older individuals have a higher risk of being overdiagnosed. Patz et al. (2014) conducted a research study on screening with low-dose computed tomography (LDCT) for lung cancer. They reported that the probability of overdiagnosis was 18.5% (95% CI, 5.4%-30.6%)

for any lung cancer detected by LDCT screening.

Marcus et al. (2006) identified 46 cases of overdiagnosis (0.75%) among 6,101 patients in the incidence follow-up of the Mayo Lung Project. This finding further supports the existence of overdiagnosis in lung cancer screening. Welch and Black (2010) investigated the risk of overdiagnosis through the detection of lung cancer using chest X-ray and/or sputum cytology. They found that approximately 51% of cases (46 in 90) may be attributed to overdiagnosis. Blom et al. (2020) analyzed different birth cohorts and their rates of overdiagnosis in screen-detected lung cancer cases. They observed that the 1950 birth cohort had a higher rate of overdiagnosis (10.5%) compared to the 1990 birth cohort (5.9%) using the cumulative excess-incidence approach. Ten Haaf and de Koning (2015) reported that 6.75% of all screen-detected cases in the chest X-ray arm and 8.62% of all screen-detected cases in the CT arm of the NLST were considered overdiagnosed. These percentages were relatively low, considering that approximately 75% of NLST participants were younger than 65, suggesting a lower potential for overdiagnosis in that particular population.

Late diagnosis might be one of the reasons for failure among patients. Early detection is necessary as lung cancer may remain incurable for patients in the advanced stage at diagnosis. Early detection trials proved a 20% reduction in lung cancer-related mortality by screening high-risk individuals with low-dose computed tomography (Vansteenkiste et al., 2012). In Wu et al. (2016), they analyzed the NLST CT scan data for lung cancer and found that the probability of true-early-detection depends more on future screening interval and the current age than on the past screening interval and the probability of true-early-detection would decrease to about 75% if the future screening interval changes from annual to biennial. Probability model to early-detection could prove major advancement as it addresses sojourn time (time duration in

preclinical states) and transition into preclinical state.

Most of the recent research dealt with overdiagnosis alone. In contrast, long term effects associated with the outcomes from true-early-detection, no-early-detection, screen-diagnosed, screen-detected cases were considered along with overdiagnosis in this research. Very few studies had dealt with probability modeling while evaluating overdiagnosis. Most of the studies relied on other characteristics, such as tumor size, cancer stage, cancer growth rate etc.. It is necessary to develop better estimates of overdiagnosis, because at the time of screening, clinicians do not know which patients have been overdiagnosed. They tend to treat all of them. Thus, overdiagnosis is associated with the problem of escalating healthcare costs. Even patients cannot benefit from unnecessary treatment, instead it is harmful. In this case, accurate estimation of the sensitivity, sojourn time, and the transition probability are very crucial. Apart from screening history and smoking status, other risk factors, such as family history, genomic aspects, inhalation of hazardous chemicals, etc. can be considered for future research.

# CHAPTER 4

# SCHEDULING THE FIRST EXAM IN LUNG CANCER SCREENING

## 4.1 Introduction

Early detection is crucial for initiating effective treatment. To improve the cure rates and increase the survival of cancer patients, screening exams should be initiated at an appropriate time. This study aims to determine the optimal timing for initiating chest x-ray for lung cancer screening in asymptomatic individuals, to facilitate early treatment. A probability method developed in Wu (2022), will be applied to the NLST chest X-ray data to identify the appropriate screening time/age; after this screening time is found, the lead time and the probability of overdiagnosis will be estimated in the (future) screening time. Although low-dose CT is the recommended modality for lung cancer screening in the United States, chest X-ray is still in use, especially in many developing countries.

Detecting lung cancer early through screening may facilitate early treatment and lead to improved long-term survival. Both chest X-rays and CT scans expose the chest to radiation, which may increase the risk of developing cancer. Therefore, it is crucial to determine the appropriate age for initiating screening and establish the frequency of recurrent screenings to minimize unnecessary exposure to radiation, harmful chemicals, and costs. A probability

method is applied to the NLST chest X-ray lung cancer screening data to identify the optimal timeline for initiating cancer screening. The optimal screening time is found by limiting the probability of incidence (from one's current age to future screening time) to a small value. The probability of incidence depends on the sensitivity, the duration of the disease-free state (transition density), and the duration of the preclinical state (sojourn time), which were previously estimated in a study in Rahman and Wu (2021).

Two other important terms in cancer screening, namely lead time and the probability of overdiagnosis, will be evaluated for individuals diagnosed with cancer during their initial screening using the probability method described in the work by Wu (2022). The time interval between detection through screening and the onset of clinical manifestations is referred to as the lead time, which represents the time duration from the detection of lung cancer through screening to the development of symptoms. In a study conducted in 2018, Liu et al. (2018) used Bayesian posterior samples of key parameters from the NLST low-dose CT data to simulate lead times by age and duration of screening intervals. On the other hand, the probability of overdiagnosis refers to the detection of disease through scheduled screening exams, but clinical symptoms would not manifest before death. It has been estimated that 18% to 67% of lung cancers detected through screening may lead to overdiagnosis, exposing many patients to unnecessary risks (Lazris and Roth, 2019). Therefore, evaluating lead time and overdiagnosis is crucial once the optimal screening time has been determined.

In the rest of the chapter, a brief review of the probability methods is presented in Section 4.2. A simulation study is conducted in Section 4.3.1 to investigate the optimal screening strategies. Bayesian inference is included in Section 4.3.2 to identify the optimal screening parameters, including screening

initiation age, lead time, and overdiagnosis. The results of the study are briefly discussed in Section 4.4.

## 4.2   Method

The contribution of this project is to apply the most current statistical method derived in Wu (2022) to accurately estimate the optimal screening time/age, lead time distribution, and overdiagnosis based on the screening frequency. We briefly review the method in Wu (2022) as follows: An asymptomatic person at current age $a_0$ has not undergone any screening yet, the first screening will occur at the age $t_0 = a_0 + t_x$, where $t_x > 0$. The objective is to find the appropriate value of $t_x$ to limit the probability of incidence to a predetermined value $p$. The goal is to restrict the probability of experiencing clinical incidence before the first screening to a small value, such as 10% or 20%. We want to ensure that there is a 90% or 80% probability of not encountering any clinical incidence before the first screening.

$$P(I_0|I_0 \cup D_0) = \frac{P(I_0)}{P(I_0 \cup D_0)} = \frac{P(I_0)}{P(I_0) + P(D_0)} = p \qquad (4.19)$$

where $P(I_0)$ is the probability of incidence in $(a_0, t_0)$ and $P(D_0)$ is the probability of detection at the first exam:

$$P(I_0) = \int_0^{a_0} w(x)[Q(a_0 - x) - Q(t_0 - x)]dx + \int_{a_0}^{t_0} w(x)[1 - Q(t_0 - x)]dx \quad (4.20)$$

$$P(D_0) = \beta \int_0^{t_0} w(x)Q(t_0 - x)dx \qquad (4.21)$$

Since $P(I_0|I_0 \cup D_0)$ is monotone increasing with $t_x$ (hence increasing with $t_0$), for any given value $p$ in (0, 1), there exists a unique solution $t_0$ such that $P(I_0|I_0 \cap D_0) = p$.

53

After $t_0$ is found, the lead time at age $t_0$, if one were diagnosed with cancer at the first exam is:

$$f_L(z|D_0) = \frac{f_L(z, D_0)}{P(D_0)}, \text{ for } z \in (0, \infty) \tag{4.22}$$

where the numerator is:

$$f_L(z, D_0) = \beta \int_0^{t_0} w(x)q(t_0 + z - x)dx \tag{4.23}$$

The probability of overdiagnosis and the probability of early detection at the first exam at one's age $t_0$ is:

$$P(OverD|D_0, T > t_0) = \int_{t_0}^{\infty} P(OverD|D_0, T = t)f_T(t|T > t_0)dt \tag{4.24}$$

$$P(TrueED|D_0, T > t_0) = \int_{t_0}^{\infty} P(TrueED|D_0, T = t)f_T(t|T > t_0)dt \tag{4.25}$$

Where the conditional PDF of a human lifetime, $f_T(t|T > t_0)$ is the same as defined in the previous chapter ((3.16)). $P(OverD|D_0, T = t)$ and $P(TrueED|D_0, T = t)$ are defined as follows:

$$P(OverD|D_0, T = t) = \frac{P(OverD, D_0|T = t)}{P(D_0|T = t)} \tag{4.26}$$

$$P(TrueED|D_0, T = t) = \frac{P(TrueED, D_0|T = t)}{P(D_0|T = t)} \tag{4.27}$$

Since, $P(D_0|T = t) = P(D_0)$, and the two numerators are derived as:

$$P(OverD, D_0|T = t) = \beta \int_0^{t_0} w(x)Q(t - x)dx \tag{4.28}$$

$$P(TrueED, D_0|T = t) = \beta \int_0^{t_0} w(x)[Q(t_0 - x) - Q(t - x)]dx \tag{4.29}$$

## 4.3 Results

To evaluate the optimal screening interval, lead time, and probability of overdiagnosis, simulation studies were initially conducted using the existing method mentioned in Section 4.2. Subsequently, a similar approach was applied to the NLST chest X-ray data for assessing these parameters.

### 4.3.1 Simulation

In the simulation study, the following scenarios were considered to estimate the optimal screening time/age, lead time, and probability of overdiagnosis:

- Four values of the probability of incidence: $p$=0.05, 0.10, 0.15, 0.20

- Three different screening sensitivities: $\beta = 0.80, 0.90, 0.95$

- Four different mean sojourn times: MST = 1.5, 2.5, 5, 10 years

- Three different current ages: $a_0$= 55, 60, 65 years

The parametric models for the transition density follow a log-Normal probability density function multiplied by 30%, while the distribution of sojourn time follows a Weibull distribution, as described in Rahman and Wu (2021). The specific forms of these models are as follows:

$$w(t|\mu, \sigma^2) = \frac{0.3}{\sqrt{2\pi}\sigma t} exp\{-(logt - \mu)^2/(2\sigma^2)\} \tag{4.30}$$

$$q(x|\alpha, \lambda) = \alpha\lambda x^{\alpha-1} exp(-\lambda x^\alpha) \tag{4.31}$$

$$Q(x) = exp(-\lambda x^\alpha) \tag{4.32}$$

Table 4.3.1: Optimal initial screening age $t_0^*$ (in years) found by binary search

| | MST=1.5 | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| p | $a_0=55$ | | | $a_0=60$ | | | $a_0=65$ | | |
| | $\beta=0.8$ | $\beta=0.9$ | $\beta=0.95$ | $\beta=0.8$ | $\beta=0.9$ | $\beta=0.95$ | $\beta=0.8$ | $\beta=0.9$ | $\beta=0.95$ |
| 0.05 | 55.08 | 55.09 | 55.09 | 60.07 | 60.08 | 60.08 | 65.07 | 65.07 | 65.08 |
| 0.10 | 55.16 | 55.19 | 55.20 | 60.15 | 60.17 | 60.18 | 65.14 | 65.16 | 65.17 |
| 0.15 | 55.26 | 55.30 | 55.32 | 60.24 | 60.27 | 60.28 | 65.22 | 65.25 | 65.26 |
| 0.20 | 55.38 | 55.43 | 55.46 | 60.34 | 60.38 | 60.41 | 65.31 | 65.35 | 65.37 |
| | MST=2.5 | | | | | | | | |
| p | $a_0=55$ | | | $a_0=60$ | | | $a_0=65$ | | |
| | $\beta=0.8$ | $\beta=0.9$ | $\beta=0.95$ | $\beta=0.8$ | $\beta=0.9$ | $\beta=0.95$ | $\beta=0.8$ | $\beta=0.9$ | $\beta=0.95$ |
| 0.05 | 55.13 | 55.15 | 55.16 | 60.12 | 60.14 | 60.15 | 65.11 | 65.13 | 65.14 |
| 0.10 | 55.29 | 55.32 | 55.34 | 60.26 | 60.29 | 60.31 | 65.24 | 65.27 | 65.29 |
| 0.15 | 55.47 | 55.53 | 55.56 | 60.42 | 60.47 | 60.50 | 65.39 | 65.43 | 65.46 |
| 0.20 | 55.68 | 55.77 | 55.82 | 60.60 | 60.68 | 60.72 | 65.55 | 65.62 | 65.65 |
| | MST=5 | | | | | | | | |
| p | $a_0=55$ | | | $a_0=60$ | | | $a_0=65$ | | |
| | $\beta=0.8$ | $\beta=0.9$ | $\beta=0.95$ | $\beta=0.8$ | $\beta=0.9$ | $\beta=0.95$ | $\beta=0.8$ | $\beta=0.9$ | $\beta=0.95$ |
| 0.05 | 55.39 | 55.45 | 55.48 | 60.32 | 60.36 | 60.38 | 65.27 | 65.30 | 65.32 |
| 0.10 | 55.89 | 56.02 | 56.08 | 60.69 | 60.79 | 60.83 | 65.57 | 65.64 | 65.68 |
| 0.15 | 56.52 | 56.75 | 56.86 | 61.14 | 61.29 | 61.37 | 65.91 | 66.03 | 66.09 |
| 0.20 | 57.32 | 57.69 | 57.88 | 61.66 | 61.89 | 62.01 | 66.31 | 66.48 | 66.56 |
| | MST=10 | | | | | | | | |
| p | $a_0=55$ | | | $a_0=60$ | | | $a_0=65$ | | |
| | $\beta=0.8$ | $\beta=0.9$ | $\beta=0.95$ | $\beta=0.8$ | $\beta=0.9$ | $\beta=0.95$ | $\beta=0.8$ | $\beta=0.9$ | $\beta=0.95$ |
| 0.05 | 55.99 | 56.14 | 56.22 | 60.79 | 60.89 | 60.95 | 65.65 | 65.73 | 65.77 |
| 0.10 | 57.47 | 57.88 | 58.09 | 61.79 | 62.05 | 62.18 | 66.41 | 66.59 | 66.69 |
| 0.15 | 59.55 | 60.31 | 60.69 | 63.05 | 63.50 | 63.73 | 67.29 | 67.59 | 67.75 |
| 0.20 | 62.05 | 63.04 | 63.52 | 64.57 | 65.22 | 65.55 | 68.31 | 68.74 | 68.96 |

In this context, $x$ represents the sojourn time in the pre-clinical state. The parameters $(\mu, \sigma^2) = (4.25, 0.015)$ were chosen such that the mode of the transition density is approximately 70 years old. For the simulation study, specific parameter values were selected to achieve the designed mean sojourn times using the Weibull distribution. The chosen values for $\alpha$ are 3.47, 1.56, 2, and 1.6, while the corresponding $\lambda$ values are 0.18, 0.202, 0.031, and 0.021. These parameter selections result in mean sojourn times of 1.5, 2.5, 5, and 10 years, respectively.

Table 4.3.1 presents the optimal initial screening age $t_0^*$ obtained using the method described in Section 4.2 and the binary search, for different values

of $p$. The analysis was conducted considering various sensitivities ($\beta$), $MST$, and current ages ($a_0$).

From Table 4.3.1, it can be observed that when $MST = 2.5$ years, the optimal screening ages under $a_0 = 55$ and $\beta = 0.9$ for different $p$ are 55.16, 55.34, 55.56, and 55.82. This indicates that if an individual seeks a 95% probability of avoiding clinical incidents before the first exam, they should undergo screening at age 55.16 (approximately two months after their current age of 55). Alternatively, if someone aims for an 80% chance of remaining free from clinical cases before the first exam, they can schedule screening at age 55.82 (approximately ten months after their current age).

Furthermore, the results demonstrate that as the screening sensitivity increases from 0.8 to 0.95, the optimal initial screening age slightly increases when other factors remain constant. However, the optimal initial screening age increases with higher incidence probability $p$ and longer mean sojourn time $MST$. The ideal first screening age $t_0^*$ is also influenced by one's current age $a_0$, with the time interval $(t_0^* - a_0)$ decreasing as $a_0$ increases, assuming other factors remain the same.

The primary objective is to find the optimal $t_0$, and after that, we can investigate the lead time distribution $f_L(z|D_0)$ and the probability of overdiagnosis $P(OverD|D_0, T > t_0^*)$ at the future screening time, $t_0^*$. In Tables 4.3.2, 4.3.3, and 4.3.4, we present the estimated mean, median, mode, and standard deviation of the lead time at the optimal first screening age $t_0^*$ for individuals with a current age of 55, 60, or 65 years, respectively.

From these tables, it is observed that the lead time distribution $f_L(z|D_0)$ is not directly influenced by the values of $\beta$, $p$, and $a_0$. However, both the lead time distribution and the probability of overdiagnosis are dependent on factors such as $t_0^*$, $w(t)$, and $Q(x)$. The findings across these three tables exhibit

Table 4.3.2: Estimated mean, median, mode and standard deviation of the lead time at optimal time $t_0^*$ when $a_0$=55

| | MST=1.5 years | | |
|---|---|---|---|
| p | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 |
| 0.05 | 0.87, 0.82, 0.67, 0.54 | 0.87, 0.82, 0.67, 0.54 | 0.87, 0.82, 0.67, 0.54 |
| 0.10 | 0.87, 0.82, 0.67, 0.54 | 0.87, 0.82, 0.67, 0.54 | 0.87, 0.82, 0.67, 0.54 |
| 0.15 | 0.87, 0.82, 0.67, 0.54 | 0.87, 0.82, 0.67, 0.54 | 0.87, 0.82, 0.67, 0.54 |
| 0.20 | 0.87, 0.82, 0.67, 0.54 | 0.87, 0.82, 0.66, 0.54 | 0.87, 0.82, 0.66, 0.54 |
| | MST=2.5 years | | |
| p | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 |
| 0.05 | 1.66, 1.59, 0.35, 1.25 | 1.66, 1.59, 0.35, 1.25 | 1.66, 1.59, 0.35, 1.25 |
| 0.10 | 1.66, 1.59, 0.35, 1.25 | 1.66, 1.59, 0.35, 1.25 | 1.66, 1.59, 0.35, 1.25 |
| 0.15 | 1.66, 1.59, 0.34, 1.25 | 1.66, 1.58, 0.34, 1.25 | 1.66, 1.58, 0.34, 1.25 |
| 0.20 | 1.66, 1.58, 0.34, 1.25 | 1.66, 1.58, 0.34, 1.25 | 1.66, 1.58, 0.33, 1.25 |
| | MST=5 years | | |
| p | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 |
| 0.05 | 1.78, 3.38, 2.25, 1.49 | 1.78, 3.38, 2.25, 1.49 | 1.78, 3.38, 2.25, 1.49 |
| 0.10 | 1.78, 3.37, 2.21, 1.49 | 1.78, 3.36, 2.20, 1.49 | 1.78, 3.36, 2.20, 1.49 |
| 0.15 | 1.78, 3.34, 2.16, 1.49 | 1.78, 3.33, 2.14, 1.50 | 1.78, 3.33, 2.13, 1.50 |
| 0.20 | 1.78, 3.31, 2.09, 1.50 | 1.77, 3.29, 2.06, 1.50 | 1.77, 3.29, 2.05, 1.51 |
| | MST=10 years | | |
| p | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 |
| 0.05 | 4.89, 5.01, 4.00, 1.62 | 4.89, 5.01, 3.98, 1.62 | 4.89, 5.01, 3.97, 1.62 |
| 0.10 | 4.89, 5.01, 3.83, 1.61 | 4.90, 5.01, 3.78, 1.61 | 4.90, 5.01, 3.76, 1.61 |
| 0.15 | 4.91, 5.01, 3.58, 1.61 | 4.92, 5.01, 3.48, 1.61 | 4.92, 5.01, 3.43, 1.61 |
| 0.20 | 4.93, 5.01, 3.24, 1.61 | 4.94, 5.01, 3.09, 1.61 | 4.94, 5.01, 3.02, 1.60 |

similar patterns, suggesting consistency in the results as follows:

i. When the MST increases, the mean, median, and mode of the lead time also increase. This indicates that a longer mean sojourn time leads to a longer expected time from the initial screening to the development of symptoms.

ii. The lead time distribution shows minimal dependence on the incidence probability $p$ and the sensitivity $\beta$ when the optimal scheduling time $t_0^*$ is used.

iii. With an increase in the current age $a_0$, the mean, median, and mode of the lead time decrease, while the standard deviation remains relatively unchanged. This suggests that as an individual's current age increases, the expected time from screening to symptom onset becomes shorter, indicating a potentially more rapid disease progression. However, the variability of lead time remains relatively consistent across different ages.

Figure 4.3.1 displays the lead time PDF curves under different factors: $p$, $\beta$, $a_0$, and $MST$. The figure consists of four panels, each showing the estimated lead time density when the optimal first screening age $t_0^*$ is used. In each panel, three factors are fixed, and only the fourth factor is allowed to vary.

The results demonstrate that, given $t_0^*$, the lead time distribution exhibits minimal changes with respect to the incidence probability $p$ and sensitivity $\beta$. However, it shows notable variation based on one's current age $a_0$ and the MST. Specifically, as $a_0$ increases, the mean, median, and mode of the lead time slightly decrease. On the other hand, as $MST$ increases, the central location of the lead time distribution shifts towards higher values.

Table 4.3.3: Estimated mean, median, mode and standard deviation of the lead time at optimal time $t_0^*$ when $a_0$=60

| p | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 |
|---|---|---|---|
| | MST=1.5 years | | |
| 0.05 | 0.85, 0.8, 0.56, 0.54 | 0.85, 0.8, 0.56, 0.54 | 0.85, 0.8, 0.56, 0.54 |
| 0.10 | 0.85, 0.8, 0.56, 0.54 | 0.85, 0.8, 0.56, 0.54 | 0.85, 0.8, 0.56, 0.54 |
| 0.15 | 0.84, 0.8, 0.56, 0.54 | 0.85, 0.8, 0.56, 0.54 | 0.85, 0.8, 0.56, 0.54 |
| 0.20 | 0.84, 0.8, 0.55, 0.54 | 0.85, 0.8, 0.55, 0.54 | 0.85, 0.8, 0.55, 0.54 |
| | MST=2.5 years | | |
| 0.05 | 1.63, 1.54, 0.21, 1.24 | 1.63, 1.54, 0.21, 1.24 | 1.63, 1.54, 0.21, 1.24 |
| 0.10 | 1.63, 1.54, 0.21, 1.24 | 1.63, 1.54, 0.21, 1.24 | 1.63, 1.54, 0.21, 1.24 |
| 0.15 | 1.63, 1.54, 0.2, 1.24 | 1.63, 1.54, 0.2, 1.24 | 1.63, 1.54, 0.2, 1.24 |
| 0.20 | 1.63, 1.53, 0.2, 1.241 | 1.62, 1.53, 0.2, 1.24 | 1.62, 1.53, 0.2, 1.24 |
| | MST=5 years | | |
| 0.05 | 1.77, 3.19, 1.82, 1.49 | 1.77, 3.19, 1.82, 1.49 | 1.77, 3.19, 1.82, 1.50 |
| 0.10 | 1.77, 3.17, 1.79, 1.50 | 1.77, 3.17, 1.78, 1.50 | 1.77, 3.17, 1.78, 1.50 |
| 0.15 | 1.77, 3.16, 1.75, 1.50 | 1.77, 3.15, 1.73, 1.51 | 1.77, 3.15, 1.72, 1.51 |
| 0.20 | 1.76, 3.14, 1.69, 1.51 | 1.77, 3.13, 1.67, 1.51 | 1.77, 3.12, 1.66, 1.51 |
| | MST=10 years | | |
| 0.05 | 4.92, 5.01, 3.42, 1.59 | 4.92, 5.01, 3.4, 1.59 | 4.92, 5.01, 3.39, 1.59 |
| 0.10 | 4.93, 5.01, 3.28, 1.59 | 4.93, 5.01, 3.24, 1.59 | 4.93, 5.01, 3.22, 1.59 |
| 0.15 | 4.94, 5.01, 3.09, 1.59 | 4.94, 5.01, 3.02, 1.59 | 4.95, 5.01, 2.99, 1.59 |
| 0.20 | 4.95, 5.01, 2.85, 1.58 | 4.96, 5.01, 2.74, 1.58 | 4.96, 5.01, 2.69, 1.58 |

Table 4.3.4: Estimated mean, median, mode and standard deviation of the lead time at optimal time $t_0^*$ when $a_0$=65

| p | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 |
|---|---|---|---|
| | MST=1.5 years | | |
| 0.05 | 0.83, 0.78, 0.42, 0.53 | 0.83, 0.78, 0.42, 0.53 | 0.83, 0.78, 0.42, 0.53 |
| 0.10 | 0.83, 0.78, 0.42, 0.53 | 0.83, 0.78, 0.42, 0.53 | 0.83, 0.78, 0.42, 0.53 |
| 0.15 | 0.83, 0.78, 0.41, 0.53 | 0.83, 0.78, 0.41, 0.53 | 0.83, 0.77, 0.41, 0.53 |
| 0.20 | 0.83, 0.77, 0.41, 0.53 | 0.83, 0.77, 0.41, 0.53 | 0.83, 0.77, 0.41, 0.53 |
| | MST=2.5 years | | |
| 0.05 | 1.59, 1.49, 0.09, 1.23 | 1.59, 1.49, 0.09, 1.23 | 1.59, 1.49, 0.09, 1.23 |
| 0.10 | 1.59, 1.49, 0.08, 1.23 | 1.59, 1.49, 0.08, 1.23 | 1.59, 1.49, 0.08, 1.23 |
| 0.15 | 1.59, 1.49, 0.08, 1.23 | 1.59, 1.49, 0.08, 1.23 | 1.59, 1.49, 0.08, 1.23 |
| 0.20 | 1.59, 1.49, 0.08, 1.23 | 1.59, 1.49, 0.08, 1.23 | 1.59, 1.49, 0.08, 1.23 |
| | MST=5 years | | |
| 0.05 | 1.76, 3.00, 1.30, 1.51 | 1.76, 3.00, 1.30, 1.51 | 1.76, 3.00, 1.30, 1.51 |
| 0.10 | 1.76, 2.99, 1.27, 1.51 | 1.76, 2.98, 1.26, 1.51 | 1.76, 2.98, 1.26, 1.51 |
| 0.15 | 1.76, 2.97, 1.23, 1.51 | 1.75, 2.97, 1.21, 1.51 | 1.75, 2.97, 1.21, 1.51 |
| 0.20 | 1.76, 2.96, 1.18, 1.51 | 1.75, 2.95, 1.16, 1.51 | 1.75, 2.95, 1.15, 1.51 |
| | MST=10 years | | |
| 0.05 | 4.96, 5.01, 2.67, 1.56 | 4.96, 5.01, 2.66, 1.56 | 4.96, 5.01, 2.65, 1.56 |
| 0.10 | 4.97, 5.01, 2.54, 1.56 | 4.97, 5.01, 2.50, 1.56 | 4.97, 5.01, 2.49, 1.56 |
| 0.15 | 4.97, 5.01, 2.38, 1.56 | 4.98, 5.01, 2.32, 1.56 | 4.98, 5.01, 2.29, 1.56 |
| 0.20 | 4.98, 5.01, 2.18, 1.56 | 4.99, 5.01, 2.1, 1.56 | 4.99, 5.01, 2.05, 1.56 |

Table 4.3.5 presents the estimated probability of overdiagnosis (in percentage) when using the optimal initial scheduling age $t_0^*$. Specifically, if an individual undergoes the first screening exam at the age $t_0^*$ provided in Table 4.3.1 and is subsequently diagnosed with cancer, the probability of overdiagnosis is given by the corresponding value in Table 4.3.5.

The probability of overdiagnosis shows an increasing trend as the MST increases. It also increases with higher incidence probability ($p$) and older age ($a_0$). However, it exhibits minimal variation with the sensitivity ($\beta$). In general, when the MST is less than or equal to one and a half years, the probability of overdiagnosis is typically less than 4%, which is considered negligible. Overall,
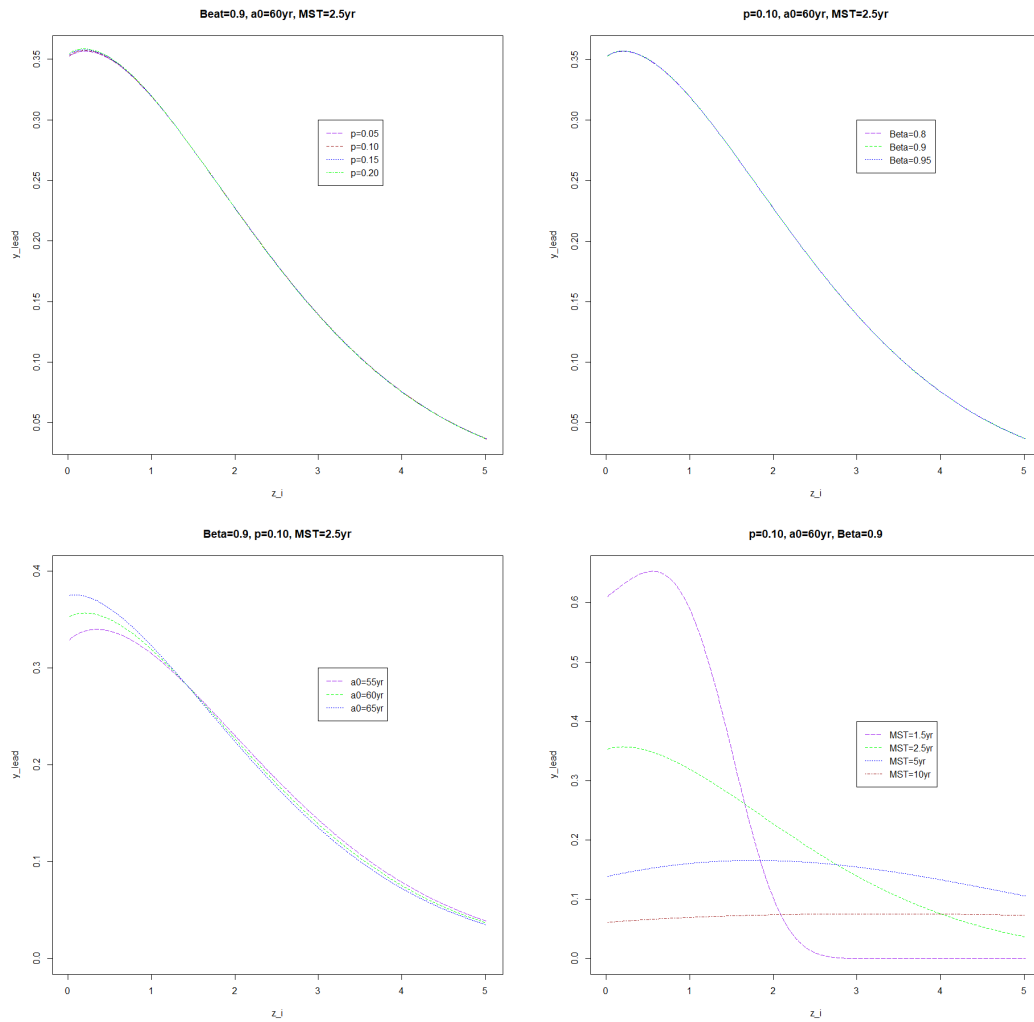
Figure 4.3.1: The PDF curves of the lead time under the four factors: changing one factor considering others as fixed

Table 4.3.5: Estimated probability of overdiagnosis (in percentage) at the initial screening age $t_0^*$

| | MST=1.5 | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| p | $a_0$=55 | | | $a_0$=60 | | | $a_0$=65 | | |
| | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 |
| 0.05 | 1.30 | 1.31 | 1.31 | 1.85 | 1.87 | 1.87 | 2.48 | 2.48 | 2.50 |
| 0.10 | 1.39 | 1.43 | 1.44 | 1.99 | 2.02 | 2.04 | 2.65 | 2.70 | 2.72 |
| 0.15 | 1.51 | 1.56 | 1.59 | 2.15 | 2.21 | 2.23 | 2.85 | 2.92 | 2.95 |
| 0.20 | 1.66 | 1.73 | 1.77 | 2.34 | 2.42 | 2.48 | 3.08 | 3.19 | 3.24 |
| | MST=2.5 | | | | | | | | |
| p | $a_0$=55 | | | $a_0$=60 | | | $a_0$=65 | | |
| | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 |
| 0.05 | 2.39 | 2.41 | 2.43 | 3.33 | 3.37 | 3.38 | 4.45 | 4.50 | 4.52 |
| 0.10 | 2.58 | 2.61 | 2.64 | 3.57 | 3.62 | 3.65 | 4.75 | 4.83 | 4.88 |
| 0.15 | 2.80 | 2.87 | 2.91 | 3.85 | 3.94 | 3.99 | 5.12 | 5.22 | 5.30 |
| 0.20 | 3.06 | 3.18 | 3.25 | 4.18 | 4.33 | 4.41 | 5.53 | 5.72 | 5.80 |
| | MST=5 | | | | | | | | |
| p | $a_0$=55 | | | $a_0$=60 | | | $a_0$=65 | | |
| | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 |
| 0.05 | 4.74 | 4.81 | 4.85 | 6.24 | 6.31 | 6.34 | 8.20 | 8.27 | 8.32 |
| 0.10 | 5.35 | 5.47 | 5.55 | 6.87 | 7.05 | 7.12 | 8.93 | 9.10 | 9.21 |
| 0.15 | 6.13 | 6.45 | 6.60 | 7.60 | 7.88 | 8.03 | 9.81 | 9.96 | 10.12 |
| 0.20 | 7.21 | 7.77 | 8.07 | 8.60 | 9.07 | 9.20 | 10.73 | 11.21 | 11.45 |
| | MST=10 | | | | | | | | |
| p | $a_0$=55 | | | $a_0$=60 | | | $a_0$=65 | | |
| | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 | $\beta$=0.8 | $\beta$=0.9 | $\beta$=0.95 |
| 0.05 | 9.39 | 9.52 | 9.65 | 10.30 | 10.53 | 10.67 | 11.41 | 11.65 | 11.77 |
| 0.10 | 10.71 | 10.46 | 10.72 | 11.04 | 11.09 | 11.11 | 12.04 | 12.18 | 12.31 |
| 0.15 | 11.46 | 11.92 | 11.76 | 12.06 | 12.16 | 12.49 | 13.81 | 13.96 | 14.53 |
| 0.20 | 12.56 | 12.84 | 12.11 | 13.07 | 13.53 | 13.61 | 14.10 | 14.34 | 14.61 |

the probability of overdiagnosis is very small, with the largest value observed being less than 15% for a MST of 10 years.

## 4.3.2   Application

To estimate the optimal screening age, lead time, and probability of overdiagnosis, 1500 Markov Chain Monte Carlo (MCMC) samples were obtained from the NLST chest X-ray data (Rahman and Wu, 2021). Equation 4.19 was applied to the NLST X-ray data for both male and female heavy smokers to determine the optimal scheduling time $t_0^*$. Once the scheduling time was estimated for the 1500 MCMC samples, the lead time distribution, probability of overdiagnosis, and true-early detection were also estimated. The values of $w(t)$, $q(x)$, and $Q(x)$ remain the same as in equations 4.30, 4.31, and 4.32, respectively. Furthermore, another crucial parameter, sensitivity was estimated using equation (2.6) described in chapter 2.

The unknown parameters $\theta = (\beta, \mu, \sigma^2, \alpha, \lambda)$ were estimated using the Markov Chain Monte Carlo (MCMC) method with a Gibbs sampler and a likelihood function. Initially, 200,000 samples were generated. After discarding the first 30,000 samples as burn-in and applying thinning every 200 iterations, a posterior sample of 500 from each chain was obtained. By running three initially overdispersed chains, a total of 1500 Bayesian posterior samples $\theta_j^*$ were obtained for each gender as the final outcome.

Table 4.3.6: Estimated initial screening age $t_0^*$ and its 95% HPD interval using the NLST X-ray data

| | MALE | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| $p$ | $a_0=55$ | | | $a_0=60$ | | | $a_0=65$ | | |
| | mean | s.e. | 95%CI | mean | s.e. | 95%CI | mean | s.e. | 95%CI |
| 0.05 | 55.042 | 0.004 | (55.035, 55.050) | 60.042 | 0.004 | (60.035, 60.050) | 65.042 | 0.004 | (65.035, 65.050) |
| 0.10 | 55.089 | 0.009 | (55.074, 55.105) | 60.089 | 0.009 | (60.074, 60.105) | 65.089 | 0.009 | (65.074, 65.105) |
| 0.15 | 55.142 | 0.014 | (55.117, 55.168) | 60.142 | 0.014 | (60.117, 60.168) | 65.142 | 0.014 | (65.117, 65.168) |
| 0.20 | 55.201 | 0.020 | (55.166, 55.237) | 60.202 | 0.020 | (60.166, 60.237) | 65.202 | 0.020 | (65.166, 65.238) |
| | FEMALE | | | | | | | | |
| $p$ | $a_0=55$ | | | $a_0=60$ | | | $a_0=65$ | | |
| | mean | s.e. | 95%CI | mean | s.e. | 95%CI | mean | s.e. | 95%CI |
| 0.05 | 55.041 | 0.004 | (55.035, 55.049) | 60.041 | 0.004 | (60.035, 60.049) | 65.041 | 0.004 | (65.035, 65.049) |
| 0.10 | 55.086 | 0.008 | (55.074, 55.103) | 60.086 | 0.008 | (60.074, 60.103) | 65.087 | 0.008 | (65.074, 65.103) |
| 0.15 | 55.137 | 0.013 | (55.118, 55.164) | 60.138 | 0.013 | (60.118, 60.164) | 65.138 | 0.013 | (65.118, 65.164) |
| 0.20 | 55.195 | 0.019 | (55.167, 55.232) | 60.195 | 0.019 | (60.167, 60.232) | 65.195 | 0.019 | (65.167, 65.232) |

In this application, a hypothetical cohort was designed for the simulation as follows: For each gender (male and female), 1500 posterior samples $\theta_j^*$, where $j = 1, 2, ..., 1500$, were obtained from the MCMC estimation using the current ages $a_0 = 55, 60, 65$. Bayesian inference was then performed to determine the optimal scheduling time for each $\theta_j^*$, based on a given incidence probability $p$. By calculating $P(I_0|I_0 \cup D_0, \theta_j^*) = p$ for each $\theta_j^*$, a corresponding scheduling age/time $t_j^*$ was determined.

The mean, standard error (s.e.), and the 95% highest posterior density (HPD) interval of the future screening age $t_j^*$ (in years) were evaluated and summarized in Table 4.3.6 using the NLST X-ray data for male and female heavy smokers. The results indicate that the optimal first screening times are very close for both genders under similar conditions, with the same current age $a_0$ and the same incidence probability $p$. However, males tend to have slightly higher optimal screening times compared to females.

After determining the optimal first screening time, the posterior distribution of the lead time was obtained as the average distribution across the pairs $(\theta_j^*, t_j^*)$, where $j = 1, 2, ..., 1500$ which has the following form:

$$f_L(z|NLST) = \frac{1}{1500} \sum_{j=1}^{1500} f_L(z|\theta_j^*)$$

The mean, median, mode, and standard deviation of the lead time, calculated using $f_L(z|NLST)$, are presented in Table 4.3.7. Generally, male heavy smokers exhibit slightly longer mean lead times compared to their female counterparts under similar conditions. Figure 4.3.2 displays the estimated lead time density curves using the NLST X-ray data, considering different current ages $(a_0)$ and incidence probabilities $(p)$. Interestingly, the lead time curves show minimal changes with respect to the incidence probability $p$ when the optimal scheduling time $t_0^*$ is employed. However, the density curves do exhibit

Table 4.3.7: Estimated mean, median, mode and standard deviation of the lead time using NLST X-ray data

| | MALE | | |
|---|---|---|---|
| p | $a_0$=55 | $a_0$=60 | $a_0$=65 |
| 0.05 | 0.68, 0.64, 0.46, 0.43 | 0.67, 0.63, 0.39, 0.43 | 0.66, 0.62, 0.31, 0.43 |
| 0.10 | 0.68, 0.64, 0.46, 0.43 | 0.67, 0.63, 0.39, 0.43 | 0.66, 0.62, 0.31, 0.43 |
| 0.15 | 0.68, 0.64, 0.46, 0.43 | 0.67, 0.63, 0.39, 0.43 | 0.66, 0.62, 0.31, 0.43 |
| 0.20 | 0.68, 0.64, 0.46, 0.43 | 0.67, 0.63, 0.39, 0.43 | 0.66, 0.62, 0.31, 0.43 |
| | FEMALE | | |
| p | $a_0$=55 | $a_0$=60 | $a_0$=65 |
| 0.05 | 0.67, 0.62, 0.35, 0.43 | 0.66, 0.61, 0.29, 0.43 | 0.66, 0.61, 0.22, 0.43 |
| 0.10 | 0.67, 0.62, 0.35, 0.43 | 0.66, 0.61, 0.29, 0.43 | 0.66, 0.61, 0.22, 0.43 |
| 0.15 | 0.67, 0.62, 0.34, 0.43 | 0.66, 0.61, 0.29, 0.43 | 0.66, 0.61, 0.22, 0.43 |
| 0.20 | 0.67, 0.62, 0.34, 0.43 | 0.66, 0.61, 0.29, 0.43 | 0.66, 0.61, 0.22, 0.43 |

variations based on the current age $a_0$: larger $a_0$ values result in higher peaks in the density curve, leading to slightly smaller mode values.

Finally, each pair $(\theta_j^*, t_j^*)$, with $j = 1, 2, ..., 1500$, is utilized to estimate the probability of overdiagnosis. The posterior mean, standard error, and 95% highest posterior density (HPD) interval of the probability (or percentage) of overdiagnosis are calculated and listed in Table 4.3.8. Additionally, the probability of true-early detection can be obtained as 1 minus the probability of overdiagnosis.

The probability of overdiagnosis at the first screening for heavy smokers, using the parameters estimated from the NLST X-ray data, is very low (less than 3%). This risk of overdiagnosis shows a slight increase with one's current age for both genders, and it is slightly higher for male heavy smokers compared to their female counterparts. Additionally, the probability of overdiagnosis slightly increases with higher values of the incidence probability $p$. It is important to note that, in this simulation, the maximum probability of overdiagnosis for both genders remains below 3%. Therefore, overdiagnosis is not a significant concern at the first screening exam using chest X-ray for heavy smokers.
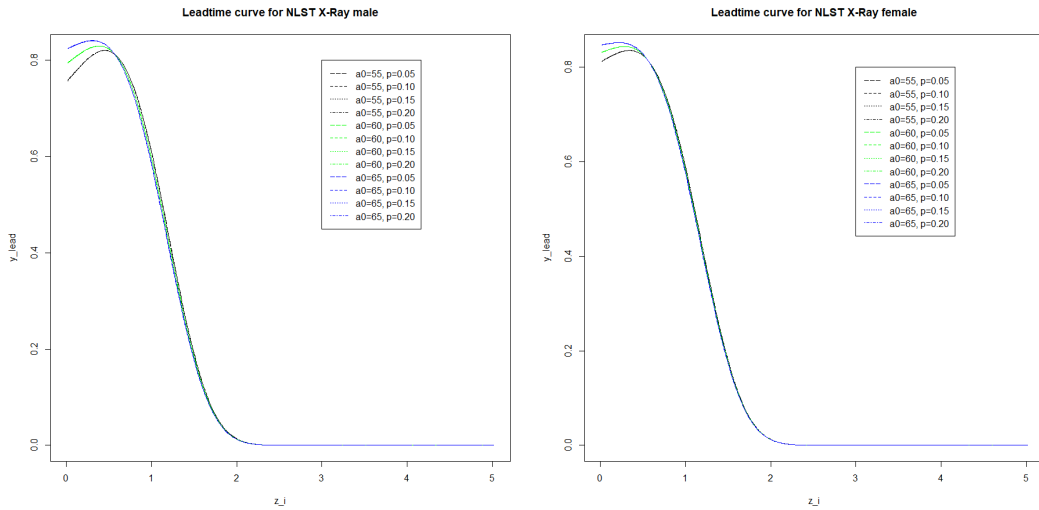
Figure 4.3.2: Lead time density curve for NLST X-ray

Table 4.3.8: Estimated mean, standard error and 95% C.I. for probability of overdiagnosis at the first exam for the NLST X-ray data (in percentage)

| | | | | MALE | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| $p$ | | $a_0=55$ | | | $a_0=60$ | | | $a_0=65$ | | |
| | mean | s.e. | 95%CI | mean | s.e. | 95%CI | mean | s.e. | 95%CI |
| 0.05 | 1.25 | 0.28 | (1.05, 1.61) | 1.78 | 0.29 | (1.06, 2.21) | 2.41 | 0.38 | (2.06, 2.97) |
| 0.10 | 1.32 | 0.25 | (1.11, 1.63) | 1.81 | 0.36 | (1.60, 2.36) | 2.52 | 0.39 | (2.16, 3.09) |
| 0.15 | 1.49 | 0.44 | (1.13, 1.65) | 1.85 | 0.38 | (1.64, 2.40) | 2.58 | 0.46 | (2.21, 3.17) |
| 0.20 | 1.53 | 0.52 | (1.23, 1.69) | 1.91 | 0.43 | (1.71, 2.44) | 2.67 | 0.58 | (2.25, 3.22) |
| | | | | FEMALE | | | | | | |
| $p$ | | $a_0=55$ | | | $a_0=60$ | | | $a_0=65$ | | |
| | mean | s.e. | 95%CI | mean | s.e. | 95%CI | mean | s.e. | 95%CI |
| 0.05 | 0.74 | 0.08 | (0.64, 0.90) | 1.03 | 0.11 | (0.89, 1.26) | 1.48 | 0.16 | (1.29, 1.82) |
| 0.10 | 0.77 | 0.08 | (0.67, 0.94) | 1.08 | 0.11 | (0.95, 1.32) | 1.55 | 0.16 | (1.35, 1.89) |
| 0.15 | 0.81 | 0.08 | (0.71, 0.98) | 1.13 | 0.11 | (0.99, 1.38) | 1.63 | 0.17 | (1.43, 1.99) |
| 0.20 | 0.85 | 0.08 | (0.75, 1.03) | 1.19 | 0.12 | (1.05, 1.45) | 1.73 | 0.17 | (1.51, 2.09) |

## 4.4 Discussion

The primary objective of this research was to determine the optimal timing for the first screening exam for asymptomatic individuals, taking into consideration their current age. The optimal first screening time for male and female heavy smokers was estimated using the NLST X-ray data, and the results were found to be consistent with the simulation study. In the simulation, it was observed that the time interval between one's current age and the first screening time slightly increases with the screening sensitivity, holding other factors constant. Additionally, it was found that the time interval increases with higher incidence probabilities. These findings align with the research conducted by Wu (2022) in her study on NLST CT scan data.

Regarding the NLST chest X-ray data used in this research, it was observed that the optimal screening age is not significantly influenced by one's current age. In other words, the difference between the optimal screening time and the current age remains relatively constant even as the current age increases. Furthermore, male heavy smokers tend to have slightly longer screening ages compared to their female counterparts. However, it is worth noting that when considering the NLST CT scan data, Wu (2022) found that female heavy smokers had longer screening ages compared to males.

In this research, it was observed that if an individual is diagnosed with cancer at the first screening exam, the lead time does not show significant changes with respect to the incidence probability and sensitivity. However, the mean, median, and mode of the lead time display a slight decrease as one's current age increases, which aligns with the findings of Wu (2022).

Furthermore, previous research by Wu et al. (2007) demonstrated that the mean lead time tends to increase as the interval between screening exams

becomes shorter. Benbassat (2021) reported a decrease in the mean lead time from 0.9 years with annual screening to 0.6 years with bi-annual screening. Regarding the comparison between male and female heavy smokers, it is worth noting that Liu et al. (2018), while analyzing NLST CT scan data, found that the mean lead time appeared longer for women than for men. However, the present study using NLST X-ray data indicates that male heavy smokers exhibit a longer mean lead time compared to their female counterparts.

The sojourn time plays a crucial role in the lead time distribution, and it is positively correlated with the mean lead time. In other words, a longer mean sojourn time corresponds to a longer mean lead time. This relationship holds true for lung cancer, as stated by Jang et al. (2013b), who noted that the distribution of the sojourn time in lung cancer is heavily skewed to the right and characterized by a large variance. Consequently, the lead time variance in lung cancer is also large, which aligns with the findings of this research.

The probability of overdiagnosis, calculated using the estimated first screening age, exhibits a positive correlation with the mean sojourn time, incidence probability, and one's current age. However, it shows only slight changes with the sensitivity, particularly when the mean sojourn time is less than 2 years. Importantly, the probability of overdiagnosis at the first screening is found to be very small. This research highlights that overdiagnosis is more closely associated with a person's lifetime, denoted as $T$. Given that the first screening occurs at a relatively younger age, it is expected to encounter small values of overdiagnosis.

# CHAPTER 5

# FUTURE WORK

This dissertation focused on estimating three essential parameters in lung cancer screening using data from the NLST chest X-ray. The estimated parameters were then utilized to infer long-term outcomes that include overdiagnosis as one outcome. The estimated parameters were also used to find the optimal age/time for screening, the distribution of lead time, and the probability of overdiagnosis at the future screening time if one would be diagnosed with cancer. In the future, I might explore working on more complicated models, such as when sensitivity depends on the sojourn time, or under other model assumptions. Additionally, I might plan to refine the likelihood function and explore alternative parametric models for these key parameters as well.

# REFERENCES

ALA (2022). American lung association. `https://www.lung.org/`.

Benbassat, J. (2021). Duration of lead time in screening for lung cancer. *BMC Pulmonary Medicine*, 21:1–8.

Blom, E. F., Ten Haaf, K., and de Koning, H. J. (2020). Trends in lung cancer risk and screening eligibility affect overdiagnosis estimates. *Lung Cancer*, 139:200–206.

Chen, Y. T., Erwin, D., and Wu, D. (2014). Over-diagnosis in lung cancer screening using the MSKC-LCSP data. *Journal of Biometrics and Biostatistics*, 5(201):2.

Chien, C. R., Lai, M. S., and Chen, T. H. (2008). Estimation of mean sojourn time for lung cancer by chest X-ray screening with a bayesian approach. *Lung Cancer*, 62(2):215–220.

Gavelli, G. and Giampalma, E. (2000). Sensitivity and specificity of chest X-ray screening for lung cancer. *Cancer*, 89(S11):2453–2456.

Hendrick, R. E. (2010). Radiation doses and cancer risks from breast imaging studies. *Radiology*, 257(1):246–253.

Jang, H., Kim, S., and Wu, D. (2013a). Bayesian lead time estimation for the Johns Hopkins lung project data. *Journal of Epidemiology and Global Health*, 3(3):157–163.

Jang, H., Kim, S., and Wu, D. (2013b). Bayesian lead time estimation for the johns hopkins lung project data. *Journal of Epidemiology and Global Health*, 3(3):157–163.

Lazris, A. and Roth, A. R. (2019). Lung cancer screening: pros and cons. *American Family Physician*, 99(12):740–742.

Liu, R., Levitt, B., Riley, T., and Wu, D. (2015). Bayesian estimation of the three key parameters in CT for the NLST data. *Journal of Biometrics and Biostatistics*, 6(5).

Liu, R., Pérez, A., and Wu, D. (2018). Estimation of lead time via low-dose CT in the National Lung Screening Trial. *Journal of Healthcare Informatics Research*, 2:353–366.

Marcus, P. M., Bergstralh, E. J., Zweig, M. H., Harris, A., Offord, K. P., and Fontana, R. S. (2006). Extended lung cancer incidence follow-up in the Mayo lung project and overdiagnosis. *Journal of the National Cancer Institute*, 98(11):748–756.

NCI (2019). Cancer stat facts: Lung and bronchus cancer. `https://www.cdc.gov/cancer/lung/index.htm`.

NIH (2020). Social security administration (SSA). `http://www.ssa.gov/OACT/STATS/table4c6.html`.

Patz, E. F., Pinsky, P., Gatsonis, C., Sicks, J. D., Kramer, B. S., Tammemägi, M. C., Chiles, C., Black, W. C., Aberle, D. R., et al. (2014). Overdiagnosis in low-dose computed tomography screening for lung cancer. *JAMA internal medicine*, 174(2):269–274.

Rahman, F. and Wu, D. (2021). Inference of sojourn time and transition density using the NLST X-ray screening data in lung cancer. *Medical research archives*, 9(5).

Tammemaegi, M. C., Church, T. R., Hocking, W. G., Silvestri, G. A., Kvale, P. A., Riley, T. L., Commins, J., and Berg, C. D. (2014). Evaluation of the lung cancer risks at which to screen ever-and never-smokers: screening rules applied to the PLCO and NLST cohorts. *PLoS medicine*, 11(12):e1001764.

Ten Haaf, K. and de Koning, H. J. (2015). Overdiagnosis in lung cancer screening: why modelling is essential. *J Epidemiol Community Health*, 69(11):1035–1039.

USPSTF (2021). Lung cancer screening. `https://www.uspreventiveservicestaskforce.org`.

Vansteenkiste, J., Dooms, C., Mascaux, C., and Nackaerts, K. (2012). Screening and early—detection of lung cancer. *Annals of Oncology*, 23:x320–x327.

Villeneuve, P. J. and Mao, Y. (1994). Lifetime probability of developing lung cancer, by smoking status, Canada. *Canadian journal of public health= Revue canadienne de sante publique*, 85(6):385–388.

Walter, S. and Day, N. (1983). Estimation of the duration of a pre-clinical disease state using screening data. *American Journal of Epidemiology*, 118(6):865–886.

Wang, D., Levitt, B., Riley, T., and Wu, D. (2017). Estimation of sojourn time and transition probability of lung cancer for smokers using the PLCO data. *J Biom Biostat*, 8(60):2.

Welch, H. G. and Black, W. C. (2010). Overdiagnosis in cancer. *Journal of the National Cancer Institute*, 102(9):605–613.

Wu, D. (2022). When to initiate cancer screening exam? *Statistics and Its Interface*, 15(4):503–514.

Wu, D., Erwin, D., and Kim, S. (2011). Projection of long-term outcomes using X-rays and pooled cytology in lung cancer screening. *Open Access Medical Statistics*, 1:13.

Wu, D., Kafadar, K., and Rosner, G. L. (2014). Inference of long term effects and overdiagnosis in periodic cancer screening. *Statistica Sinica*, pages 815–831.

Wu, D., Kafadar, K., Rosner, G. L., and Broemeling, L. D. (2012). The lead time distribution when lifetime is subject to competing risks in cancer screening. *The International Journal of Biostatistics*, 8(1).

Wu, D. and Kim, S. (2020). Problems in the estimation of the key parameters using MLE in lung cancer screening. *Journal of clinical research and reports*, 5(3).

Wu, D., Liu, R., Levitt, B., Riley, T., and Baumgartner, K. (2016). Evaluating long-term outcomes via Computed Tomography in lung cancer screening. *J Biom Biostat*, 7(313):2.

Wu, D., Rosner, G. L., and Broemeling, L. (2005). MLE and bayesian inference of age-dependent sensitivity and transition probability in periodic screening. *Biometrics*, 61(4):1056–1063.

Wu, D., Rosner, G. L., and Broemeling, L. D. (2007). Bayesian inference for the lead time in periodic cancer screening. *Biometrics*, 63(3):873–880.

Zelen, M. and Feinleib, M. (1969). On the theory of screening for chronic diseases. *Biometrika*, 56(3):601–614.

# CURRICULUM VITA

**NAME:**           Farhin Rahman

**ADDRESS:**      Department of Bioinformatics and Biostatistics

University of Louisville, Louisville, KY 40292

**EDUCATION:**

Bachelor of Science in Applied Statistics,

University of Dhaka, Dhaka, Bangladesh, 2011

Master of Science in Statistics,

Ball State University, Indiana, USA, 2017

**PUBLICATIONS:**

**Rahman, F.**, Wu, D. (2021) Inference of

Sojourn Time and Transition Density using the

NLST X-ray Screening Data in Lung Cancer.

*Medical research archives.*

**Rahman, F.**, Begum, M. (2018)

Survival Analysis of Recurrent Events on

Prostate Cancer: Facts from Cancer Genome.

*Journal of Statistical Research.*

**PRESENTATIONS:**

The Statistics and data Science Conference,

*April 2022.* Inference of Sojourn Time and

Transition Density using the NLST X-ray

Screening Data in Lung Cancer.

Summer Public Health Workshop, *July 2021*.

Estimation in Lung Cancer Screening.

ASA-KY Chapter Meeting, *April 2021*.

Inference of Sojourn Time and Transition

Density using the NLST X-ray Screening

Data in Lung Cancer.

Department of Bioinformatics and Biostatistics

Spring Seminar Series, University of Louisville,

*March 2021.* Inference of Key Parameters using

the NLST X-ray Screening Data in Lung Cancer.

**HONORS**

**AND AWARDS:**

*Graduate Student Assistantship*, Resources

of Academic Achievements, University of

Louisville, August 2018 - December 2019

*Research Assistantship*, Department of

Bioinformatics and Biostatistics, University of

Louisville, January 2020 - August 2023

*Graduate Assistantship*, Department of Math-

ematical Sciences, Ball State University, August

2015 - May 2017

*Research Assistantship*, Center for Business

and Economic Research, Ball State University,

August 2017 - June 2018