

NBER WORKING PAPER SERIES

BOUNDS IN COMPETING RISKS MODELS AND THE WAR ON CANCER

Bo E. Honoré
Adriana Lleras-Muney

Working Paper 10963
<http://www.nber.org/papers/w10963>

NATIONAL BUREAU OF ECONOMIC RESEARCH
1050 Massachusetts Avenue
Cambridge, MA 02138
December 2004

We would like to thank Jaap Abbring, Josh Angrist, Eric J. Feuer, Marco Manacorda, researchers at the National Cancer Institute, and seminar participants at CAM at the University of Copenhagen, London School of Economics, Massachusetts Institute of Technology, Princeton University, University College London, and the Harvard-MIT-Boston University Health seminar for their suggestions. This research was supported by the National Institute on Aging, Grant Number K12-AG00983 to the National Bureau of Economic Research (Adriana Lleras-Muney) and by the National Science Foundation, The Gregory C. Chow Econometric Research Program at Princeton University, and the Danish National Research Foundation, through CAM at The University of Copenhagen (Bo Honoré). The views expressed herein are those of the author(s) and do not necessarily reflect the views of the National Bureau of Economic Research.

© 2004 by Bo E. Honoré and Adriana Lleras-Muney. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

Bounds in Competing Risks Models and the War on Cancer

Bo E. Honoré and Adriana Lleras-Muney

NBER Working Paper No. 10963

December 2004

JEL No. I10, C40

ABSTRACT

In 1971 President Nixon declared war on cancer and increased the federal funds allocated to cancer research dramatically. Thirty years later, many have declared this war a failure. Overall cancer statistics confirm this view: age-adjusted mortality in 2000 was essentially unchanged from the early 1970s. At the same time, age-adjusted mortality rates from cardiovascular disease have fallen quite dramatically. Since the causes underlying cancer and cardiovascular disease are likely to be correlated, the decline in mortality rates from cardiovascular disease may be somewhat responsible for the rise in cancer mortality. It is natural to model mortality with more than one cause of death as a competing risks model. Such models are fundamentally unidentified, and it is therefore difficult to get a clear picture of the progress in cancer. This paper derives bounds for aspects of the underlying distributions under a number of different assumptions. Most importantly, we do not assume that the underlying risks are independent, and impose weak parametric assumptions in order to obtain identification. The theoretical contribution of the paper is to provide a framework to estimate competing risk models with interval data and discrete explanatory variables, both of which are common in empirical applications. We use our method to estimate changes in cancer and cardiovascular mortality since 1970. The estimated bounds for the effect of time on the duration until death for either cause are fairly tight and we find that trends in cancer show much larger improvements than previously estimated. For example, we find that time until death from cancer increased by about 10% for white males and 20% for white women.

Bo E. Honoré
Department of Economics
Princeton University
Princeton, NJ 08544-1021
honore@princeton.edu

Adriana Lleras-Muney
Department of Economics
Princeton University
Princeton, NJ 08544-1021
and NBER
alleras@princeton.edu

1 Introduction

In 1971 President Nixon declared war on cancer. As a result the Nixon administration created a National Cancer Program administered by the National Cancer Institute, and increased the federal funds allocated to cancer research dramatically.¹ Thirty years later, however, many have declared this war a failure (Bailar and Smith (1986), Bailar and Gornik (1997), etc). Overall cancer statistics confirm this view. Age-adjusted mortality from cancer increased from 198.7 (per 100,000) in 1973 to 213 in 1993, and then it fell to about its 1973 levels (198.6) in 2000. Incidence rates show a similar pattern, increasing from 385 in 1973 to 509.85 in 1992, and then decreasing to 477 in 2000 (SEER (2004)).

At the same time, age-adjusted mortality rates from cardiovascular disease have fallen quite dramatically. (See Figure 3.) It has been hypothesized that the decline in mortality rates from cardiovascular disease is somewhat responsible for the rise in cancer mortality. In other words, perhaps if there had been no progress in cardiovascular disease, we might have observed different trends in cancer mortality. The intuition behind the hypothesis that observed cancer trends are biased is that the fall in mortality rates from cardiovascular disease leaves more and perhaps different individuals at risk for cancer. Indeed for younger individuals, for whom cardiovascular disease is not a large competing risk, there have been large improvements in cancer: since 1973, cancer mortality for children and adolescents (under age 20) has fallen by more than 50% across all types of cancers, and it fell by 20% for young adults ages 20 to 44. Moreover these reductions have occurred in spite of the increases in cancer incidence for both groups (Doll (1991)). The same is not true for older adults. Although it has long been recognized that dependent competing risks can affect trends in cancer mortality, no estimates of cancer trends exist that account for this possibility.² In fact in 1990, the Extramural Committee to Assess Measures of Progress Against Cancer recommended “additional research on how cancer statistics are affected by changes in other causes of death.”

This paper derives bounds for aspects of the underlying distributions under a number of different assumptions. Most importantly, we do not assume that the underlying risks are independent, and impose very weak parametric assumptions in order to obtain identification. The theoretical contribution of the paper is to provide a framework to estimate competing risk models with interval

¹The National Cancer Institute’s budget is approximately \$4.3 billion (or 18% of the budget for the NIH).

²Chiang (1991), Rothenberg (1994) and Llorca and Delgado-Rodriguez (2001) have investigated the effects of cardiovascular mortality trends on trends in cancer mortality. However as Wohlfart and Andersen (2001) point out, these authors assume that risks are independent in their analyses.

data and discrete explanatory variables, both of which are common in empirical applications. There are a number of economic applications of the competing risks model in economics. For example, Flinn and Heckman (1982) investigated the duration of unemployment where an employed individual could terminate a spell of unemployment either by finding a job or by leaving the labor market. Katz and Meyer (1990) used the competing risks model to study the probability of leaving unemployment through recalls and new jobs. Other applications include studies of age at marriage or cohabitation (Berrington and Diamond (2000)), Ph. D. completion (Booth and Satchell (1995)), and mortgage termination (Deng, Quigley, and Van Order (2000)). The competing risks model is also closely related to the Roy (1951) model studied in Heckman and Honoré (1990) and Heckman, Smith, and Clements (1997).

This framework is then applied to mortality data from the US to estimate the trends in cancer mortality, which are the most widely used measure of overall progress against cancer.³ We find that trends in cancer show much larger improvements than previously estimated.

2 Data

We use mortality rates by single year of age, gender, race (black and white) and cause of death. These were calculated by matching population data from the Census Bureau and number of deaths from the Multiple Cause of Death Mortality files from 1970, 1980, 1990 and 2000. We computed mortality rates for three causes of death: cardiovascular disease (hereafter CVD), cancer and all other causes. (For data sources and details see the appendix.) We restrict the sample to individuals over age 45, so all the results we present are conditional on survival to that age. For 1970, population counts exist by single year of age up to age 79, and by 5-year intervals over age 80. To obtain consistent results over time, we therefore censor durations for all years at age 80.

Table 1 presents summary statistics of the data (prior to censoring at age 80) for each census year and for four demographic groups defined by gender and race. It documents the well-known

³There are several measures used to assess progress in cancer, including age-adjusted incidence rates, 5-year survival rates conditional on diagnosis, and mortality rates. Both survival rate conditional on diagnosis and incidence rates are affected by improvement in diagnosis technology. Better diagnostic tools allow for detection of tumors at earlier stages, generating a mechanical increase in survival rates that does not reflect improvements in prevention or treatment (Welch, Schwartz, and Woloshin (2000)). Similarly, improved detection increases observed incidence, even though disease rates may not have changed. Additionally, diagnosis is a function of access to care, further complicating the interpretation of changes in incidence and 5-year survival rates. For these reasons, when reporting to the Senate Appropriations committee in 1990, the Extramural Committee to Assess Measures of Progress against Cancer concluded that age-specific cancer mortality is the best measure of progress against cancer.

patterns in mortality. As of 1970, between 55 and 70% of individuals died from CVD. However there were large differences across demographic groups in age at death from all causes and from cancer and CVD: white women lived the longest, followed by white men, black women and lastly black men. From 1970 to 2000, all groups experienced an increase in the age at death; and the share of individuals dying from cardiovascular disease fell dramatically while the share dying from cancer increased for all groups (although it fell in the 1990s for all except white men). But again there are some important differences across groups: the increase in life expectancy was largest for black females, the reductions in the percentage of CVD deaths were largest for whites and the percentage increases in deaths from cancer were largest for black men. Because of these differences we analyze the results separately for each group.

With our data we can calculate the observed hazard rates using a discrete time Kaplan–Meier estimator. Figures 5 and 6 show these sub-hazards for white males, white females, black males and black females, for cancer and CVD separately. These hazard rates present in more detail the same trends that the summary statistics show. Hazard rates from CVD declined quite significantly in every decade for all groups. On the other hand, there is no discernible trend in cancer hazard rates. It is also clear that hazard rates are fairly different across demographic groups. From these graphs we also note that hazard rates are much more volatile among blacks, especially at older ages. This is true for both causes of death, but it is more pronounced for cancer rates. Censoring at age 80 alleviates the problem somewhat since hazard rates become even more volatile for older ages (not shown).

3 Competing Risks

In this section we review the theory on competing risks, illustrating issues and methods in the context of cancer and cardiovascular mortality and using the data we just described.

3.1 Set-up

Formally, a competing risks model is a duration model where the observed duration is the shortest of a number of latent durations. In addition it is typically also assumed that the identity of the shortest duration is observed. Mathematically, we observe T and δ where

$$(T, \delta) = (\min \{T_1, T_2, \dots, T_K\}, \arg \min \{T_1, T_2, \dots, T_K\}).$$

See, for example, Kalbfleisch and Prentice (1980) or Crowder (2001). Much of the terminology in this literature is motivated by medical applications where T_k could be the unobserved (latent)

duration until death from a specific cause (risk) such as cancer or cardiovascular disease, T the observed duration until death and δ the cause of death. In order to simplify the exposition and to present the theory related to the specific case we analyze, we will focus on the case where $K = 2$ in what follows. The general case requires no additional ideas, but the notation is substantially more cumbersome in that case.

In this paper we will use the notation

$$T^* = \min \{T_1, T_2\}, \quad \delta = 1 \{T_1 < T_2\}$$

and the object of interest will be features of the distribution of (T_1, T_2) given a set of explanatory variables X . Knowledge of the joint distribution of the unobserved, latent distributions T_1 and T_2 (given X) allows one to answer policy questions that one could not answer on the basis of the distribution of (T^*, δ) (given X). For example, the latter will not allow one to evaluate the effect of eliminating one of the risks on the distribution of the duration until death.

As discussed below, applications of competing risks models have often, though not always, assumed that the underlying latent durations are statistically independent. While such an assumption is reasonable in some contexts, there are at least two related reasons why one could suspect it to be violated in specific situations.

The first reason why the latent durations might be dependent is that the same underlying process affects both risks. In the case of CVD and cancer, there are several common risk factors that affect both. The American Heart Association lists smoking, drinking alcohol in large amounts, and obesity as factors that increase the likelihood of coronary heart disease, stroke, high blood pressure and hypertension. Moderate alcohol consumption and exercise on the other hand reduce blood pressure and coronary heart disease. The National Cancer Institute and the American Cancer Society also document that the same factors affect the risk of certain cancers. Smoking increases cancers of the respiratory system, as well as other cancers. Obesity increases the risk of cancer of the uterus, breast and prostate cancer among others. Excessive alcohol use increases the risk of cancer of the mouth, pharynx, larynx, esophagus, liver, and breast. Exercise is thought to reduce the risk of colon and breast cancers, and moderate alcohol consumption may lower the risk of leukemia, skin, breast and prostate cancers. This evidence suggests that at the individual level, cancer and CVD are not independent risks.

Additionally, heterogeneity across individuals can cause the underlying latent durations to be dependent even if the risks are independent for every individual in that population (Vaupel and Yashin (1999)). There is substantial evidence of genetic differences across individuals with respect to their susceptibility to both CVD (Nabel (2003)) and cancer (e.g. Lynch and de la Chapelle (2003)),

Wooster and Weber (2003)).⁴ This will cause the latent duration until death from CVD and cancer to be correlated. Furthermore there are large differences in the population in terms of exposure to environmental factors and behaviors that increase particular death risks. For example in 2000, high school dropouts were more than twice as likely to smoke than college-educated individuals; women below poverty level were twice as likely as women in the highest income levels to be obese; married individuals were less likely to exercise than those who have never married; and Hispanics were less likely than non-Hispanics to drink (Schoenborn, Adams, Barnes, Vickerie, and Schiller (2004)).

This suggests that it is interesting to consider competing risks models with dependent latent durations.

3.2 Identification

The identification of the competing risks model is tricky. The key result in this literature is that for any joint distribution of (T_1, T_2) , there exists (unique) univariate distribution for S_1 and S_2 , such that if S_1 and S_2 are independent, then the distribution of $(\min\{T_1, T_2\}, 1\{T_1 < T_2\})$ equals that of $(\min\{S_1, S_2\}, 1\{S_1 < S_2\})$ See Cox (1962) and Tsiatis (1975). In other words, for every dependent distribution of (T_1, T_2) , one can find an independent distribution that generates observationally equivalent data. Since this exercise can be carried out conditional on a set of explanatory variables X , the relationship between T_1 and T_2 conditionally on X is fundamentally unidentified, and it is not possible to use observational data only to test whether or not the risks are dependent. It is therefore necessary to make additional assumptions if one wants to answer questions that require exact knowledge of the joint distribution of (T_1, T_2) .

Broadly speaking, there have been three approaches to dealing with the identification problem in competing risks. The first is to make no additional assumptions and to estimate bounds for the object of interest, for example the marginal distributions of the underlying durations. The second approach is to assume that the risks are independent (conditional on a set of observed covariates) in which case estimation of competing risks models amounts to estimation of duration models with random censoring. The third broad approach is to specify a parametric or semiparametric model for (T_1, T_2) conditional on the covariates. The approach taken in this paper is a combination of the first and the third approach.

If one is willing to assume independence then it is straightforward to estimate the hazard function for each of the underlying distributions. For the case of cancer, the hazard rates in Figure

⁴See the web pages of the American Heart Association and the National Cancer Institute for additional cites.

5 are sufficient to conclude that there has been a very small improvement in cancer mortality, if any at all. Of course, imposing independence when the risks are indeed dependent, will result in inconsistent estimates of the cause-specific hazard rates and of the effect of covariates on those hazards.⁵ Given that the medical evidence suggests that CVD and cancer are dependent, it is therefore not possible to reach definite conclusions by looking at the observed hazards, as we did above.

Alternatively, one can make no assumptions on the joint distribution of the underlying durations, and estimate bounds on the objects of interest. Following the approach of, for example, Peterson (1976) and Manski (2003), it is straightforward to generate bounds on the marginal distributions of T_1 and T_2 . These bounds are given in Peterson (1976), who also provides bounds on the joint distribution of T_1 and T_2 . It is easy to understand the basic idea behind these bounds. For example suppose that by age 60, 15% of individuals have died of CVD and 10% have died of cancer. The survival rate⁶ from cancer at age 60 can be bounded between 75 and 90%. Although this approach is very appealing, the nonparametric bounds are generally very wide (see the numerical example in Peterson (1976)), making it difficult to draw conclusions. In Figure 4 we present the bounds for the survival from cancer in 1970 and 2000 for our four demographic groups. It is evident from these graphs that it is not possible to make any statement about whether survival from cancer increased or decreased in this period.

The results presented in Figure 4 and the potential problems with assuming independence, suggest that it might be fruitful to ask what features of the conditional distribution of (T_1, T_2) , given some explanatory variable X , can be identified if one is willing to impose restrictions on those conditional distributions. At the extreme, one could specify a fully parametric model and estimate the parameters of such a model by maximum likelihood. This is the approach taken in most of the applications cited in the introduction. The weakness of a fully parametric approach is that the results may be entirely driven by the functional form assumptions. A number of papers have therefore studied identifiability of semiparametric competing risks models.

Heckman and Honoré (1989) show (essentially) that with a mixed proportional hazard model

⁵For example, when studying mortality by cause, one may be willing to assume that drug X affects only S_1 but by imposing independence we will estimate that drug X also affects S_2 . Slud and Byar (1988) provide such an example. Vaupel and Yashin (1999) illustrate the problems that arise if one assumes independence in the presence of unobserved population heterogeneity (which results in dependent population hazards).

⁶The epidemiology literature on cancer often uses the term “survival rate” to refer to the fraction of people who are alive five years after being diagnosed. In this paper we do not condition on diagnosis and we do not only consider five year periods.

or an accelerated failure time model on the marginal distributions of T_1 and T_2 , the full model is identified if one is willing to assume that the support of the effect of X on the hazard functions for T_1 and T_2 is \mathbb{R}_+^2 . A recent paper by Abbring and van den Berg (2003) relaxes these conditions somewhat by showing that the unbounded support assumption can be dispensed with if one is willing to make additional assumptions. However, as discussed by Crowder (2001) the conditions for identification are restrictive and often unrealistic as the covariates of interest have bounded support and are not continuous in many applications. For example, analyses of mortality use data from death certificates, which contain demographic information that is all categorical, such as race, gender and marital status. Moreover, the proofs in Heckman and Honoré (1989) and Abbring and van den Berg (2003) rely crucially on the duration, T , being observed exactly. However, the durations are observed in groups in many data sets. This raises the question of what can be learned in competing risks models if one is willing to impose restrictions that are weaker than those in Heckman and Honoré (1989) and Abbring and van den Berg (2003). This is the subject of the next section.

Competing risks models are a subset of sample selection models. The research presented here is therefore closely related to the literature on bounds in sample selection models (see for example Manski (1990)), although the results here take advantage of the special structure of the competing risks model.

4 Bounds in Some Specific Competing Risks Models

As mentioned above, one of the motivations for this paper is that many data sources contain interval observations of durations, whereas the results on identification of semiparametric competing risks models assume that durations are observed exactly. Following, for example, Prentice and Gloeckler (1978) and Meyer (1990), we assume that (T_1, T_2) has a continuous positive density conditional on X , but that $T^* = \min\{T_1, T_2\}$ is grouped so we observe events like (T, δ, X) , where $T = t_k$ if $t_k < T^* \leq t_{k+1}$ for $k = 1, \dots, M$ and $t_{M+1} = \infty$. In the following we assume M is finite, so that there is only a finite number of possible outcomes. We also assume that δ is unobserved when $T^* > t_M$. In other words, we allow T^* to be censored at t_M .

The main methodological contribution of the research presented in this section is to show how parametric assumptions can help tighten the bounds on the object of interest in unidentified competing risks models. This is interesting because the nonparametric bounds that make no assumptions can be quite wide. Since different assumptions will lead to different sets of identified regions, we will consider a number of examples. In each of the examples, we will use the fact that for any

distribution of (T_1, T_2) given X , there exist an observationally equivalent discrete distribution for which the probability of a tie is 0. This follows from the fact that only a discretized version of T is observed. If X can take a finite number of values, this means that for all the cases we consider, there will be an observationally equivalent case in which the vector of all the random variables has a discrete distribution with a finite number of points of support.

4.1 The effect of explanatory variables with parametric restrictions.

We first consider the case where a binary explanatory variable, X , has a multiplicative effect on both of the latent distributions,

$$(T^*, I) = \begin{cases} (\min\{S_1, S_2\}, 1\{S_1 < S_2\}) & \text{for } X = 0, \\ (\min\{\alpha S_1, \beta S_2\}, 1\{\alpha S_1 < \beta S_2\}) & \text{for } X = 1, \end{cases} \quad (1)$$

where (S_1, S_2) is independent of X , and the multiplicative effect, α , is the main object of interest. In the next section we also consider the case where no assumption is made on the effect of X on T_2 . This model is an example of an accelerated failure time model, which is commonly used to describe mortality. It was originally introduced by Cox (1972), who gave it a physical interpretation in the context of mortality. From the equivalence between proportional hazard models with Weibull baseline hazards and Weibull accelerated failure time models, it follows that a model where the marginals obey a Weibull proportional hazard assumption, are consistent with our functional form assumption. It is also a special case of the kind of general sample selection models that have been considered in the econometric literature. Specifically, if the durations are not grouped, then one can write the model in (1) as a switching regression model. See Amemiya (1985). Specifically, let $\varepsilon_k = \log(S_k)$ and consider $\log(T_1)$

$$\log(T_1) = X \cdot \log(\alpha) + \varepsilon_1$$

where $\log(T_1)$ is observed only if

$$X \cdot (\log(\beta) - \log(\alpha)) + (\varepsilon_2 - \varepsilon_1) < 0$$

The standard sufficient conditions for identification of such models require that X has “full rank” conditional on the probability that the selection criterion is satisfied (i.e. conditional on the so-called propensity score). See for example Ahn and Powell (1993). This sufficient condition is not satisfied here. Moreover, it is clear that a model with a finite number of points of support for the explanatory variable and a discrete outcome variable will not be point-identified (by the same

intuition why a semiparametric discrete choice model is not identified if the explanatory variables take only a finite number of values).

Because the parameters in (1) are not point-identified, we will construct bounds on them. To see that one can obtain bounds on the parameters under (1), consider a simple case in which observations are censored after 2 time-periods, and one observes

$$\begin{aligned} P(T = 0, I = 0 | X = 0) &= P(T = 1, I = 0 | X = 0) = \\ P(T = 0, I = 1 | X = 0) &= P(T = 1, I = 1 | X = 0) = \\ P(T = 2 | X = 0) &= \frac{1}{5} \end{aligned}$$

If $a = b = 2$ then this has a number of implications for the distribution of (T, I) when $X = 1$. For example, all observations that were censored when $X = 0$ will still be censored when $X = 1$, as will observations with $T = 1$ (so T^* is between 1 and 2). On the other hand, none of the observations with $T = 0$ (so T^* is between 0 and 1) will be censored. This means that the probability of censoring must be 0.6 when $X = 1$. So, if one observes that the probability of censoring when $X = 1$ is 0.5, then one would rule out $a = b = 2$. Of course, in this case there are additional constraints, for example $P(T = 0 \text{ or } 1, I = 0 | X = 1) = 0.2$. The main insight in this section is to keep track of all such implications for a given (a, b) . To do this, we make use of the fact that for any parameter value which is consistent with the observed distribution of the data, there is a discrete distribution of the underlying random variables that makes it consistent with the data. In asking whether particular values of α and β are consistent with the observed distribution of the data, there is therefore no loss in generality by assuming that the underlying distributions are discrete (with support that depends on α and β). The points of support will be denoted by (s_1, s_2) , and the associated probabilities by $p(s_1, s_2)$. In this case, the relevant probabilities are

$$P(t < S_1 < t + 1, S_1 < S_2) \tag{2}$$

$$P(t < S_2 < t + 1, S_2 < S_1) \tag{3}$$

(corresponding to $X = 0$) and

$$P(t < \alpha S_1 < t + 1, \alpha S_1 < \beta S_2) = P\left(\frac{t}{\alpha} < S_1 < \frac{t+1}{\alpha}, S_1 < \frac{\beta}{\alpha} S_2\right) \tag{4}$$

$$P(t < \beta S_2 < t + 1, \beta S_2 < \alpha S_1) = P\left(\frac{t}{\beta} < S_2 < \frac{t+1}{\beta}, S_2 < \frac{\alpha}{\beta} S_1\right) \tag{5}$$

(corresponding to $X = 1$).

In order to construct the relevant points of support, consider the set of numbers $\{0, 1, 2, 3, \dots, t_{Max}\} \cup \{0, \alpha^{-1}, 2\alpha^{-1}, 3\alpha^{-1}, \dots, t_{Max}\alpha^{-1}\}$. Label this set $\{q_1, q_2, \dots, q_K\}$. These are the relevant numbers as far as the marginal distribution of T_1 is concerned. Also consider the set of numbers

$\{0, 1, 2, 3, \dots, t_{Max}\} \cup \{0, \beta^{-1}, 2\beta^{-1}, 3\beta^{-1}, \dots, t_{Max}\beta^{-1}\}$. Label this set $\{r_1, r_2, \dots, r_L\}$. These are the relevant numbers for the marginal distribution of T_2 .

The first two graphs in Figure 1 depict the events in equations (2) and (3), and in equations (4) and (5), respectively. The dashed lines in the graphs correspond to the numbers $\{0, 1, 2, 3, \dots, t_{Max}\}$ and the dotted lines to $\{0, \alpha^{-1}, 2\alpha^{-1}, 3\alpha^{-1}, \dots, t_{Max}\alpha^{-1}\}$ and $\{0, \beta^{-1}, 2\beta^{-1}, 3\beta^{-1}, \dots, t_{Max}\beta^{-1}\}$.

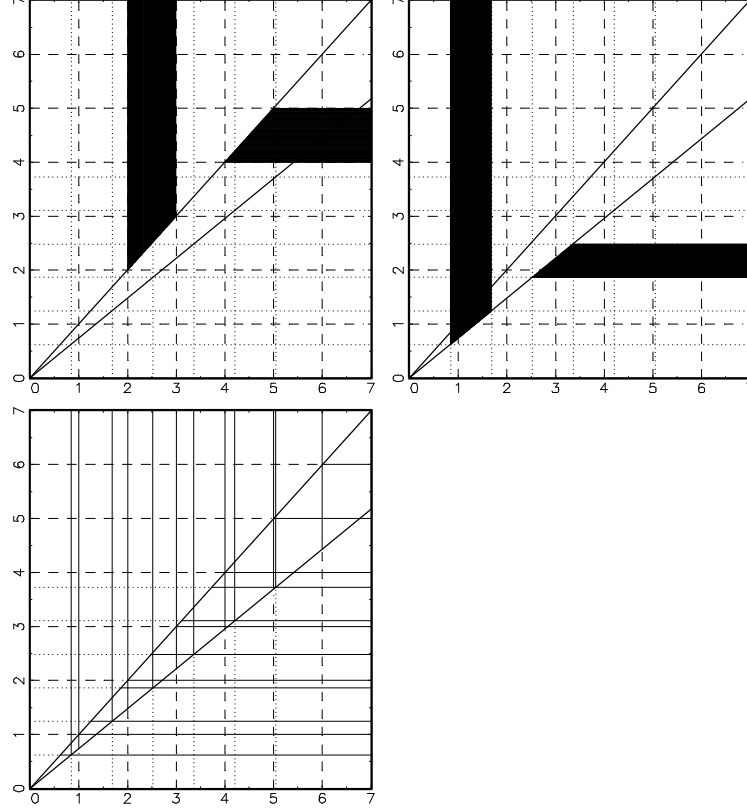


Figure 1: Illustration of Points of Support for Case (a)

It is clear that the probabilities of those events would be unchanged if one redistributed probability within each of the polygons depicted (in solid lines) in the third graph. There is therefore no loss of generality in assuming that the distribution of (S_1, S_2) is discrete, with one point of support in each of the regions.

The identified region for (α, β) is the set of (a, b) such that there exists $p(s_1, s_2)$ satisfying⁷

$$\sum_{\substack{t_k < s_1 < t_{k+1} \\ s_2 > s_1}} p(s_1, s_2) = P(T = t_k, I = 1 | X = 0), \quad (6)$$

$$\sum_{\substack{t_k < s_2 < t_{k+1} \\ s_1 > s_2}} p(s_1, s_2) = P(T = t_k, I = 0 | X = 0), \quad (7)$$

$$\sum_{\substack{t_k < a s_1 < t_{k+1} \\ b s_2 > a s_1}} p(s_1, s_2) = P(T = t_k, I = 1 | X = 1), \quad (8)$$

$$\sum_{\substack{t_k < b s_2 < t_{k+1} \\ a s_1 > b s_2}} p(s_1, s_2) = P(T = t_k, I = 0 | X = 1), \quad (9)$$

$$\sum_{s_1, s_2} p(s_1, s_2) = 1, \quad p(s_1, s_2) \geq 0 \quad (10)$$

(where the first four equations hold for all $k = 1, \dots, M$).

These equations have exactly the same structure as the constraints of a linear programming problem. Analogous to Honoré and Tamer (2004) and Molinari (2004), one can check whether a feasible solution to such a linear programming problem exists for a given a and b by solving an auxiliary linear programming problem and checking whether its optimal value is 0 (the alternative being that it is negative). We will show that as suggested in Honoré and Tamer (2004), one can consistently estimate the identified region for (α, β) by maximizing the optimal value in the sample analogs to the auxiliary linear programming problem.

Specifically, for given a and b consider the linear programming problem

$$f(a, b) = \max_{\{v_i\}, \{p(\cdot, \cdot)\}} \sum -v_i \quad (11)$$

subject to

$$\begin{aligned} v_k + \sum_{\substack{t_k < s_1 < t_{k+1} \\ s_2 > s_1}} p(s_1, s_2) &= P(T = t_k, I = 1 | X = 0) & k = 1, \dots, M, \\ v_{M+k} + \sum_{\substack{t_k < s_2 < t_{k+1} \\ s_1 > s_2}} p(s_1, s_2) &= P(T = t_k, I = 0 | X = 0) & k = 1, \dots, M, \end{aligned}$$

⁷Imposing $\beta = 1$ in this example, will give the identified region for α , under the exclusion restriction that X has no effect on T_2 .

$$\begin{aligned}
v_{2M+k} + \sum_{\substack{t_k < as_1 < t_{k+1} \\ bs_2 > as_1}} p(s_1, s_2) &= P(T = t_k, I = 1 | X = 1) \quad k = 1, \dots, M, \\
v_{3M+k} + \sum_{\substack{t_k < bs_2 < t_{k+1} \\ as_1 > bs_2}} p(s_1, s_2) &= P(T = t_k, I = 0 | X = 1) \quad k = 1, \dots, M, \\
v_{4M+1} + \sum_{s_1, s_2} p(s_1, s_2) &= 1, \quad p(s_1, s_2) \geq 0 \quad \text{for all } (s_1, s_2), \\
v_i &\geq 0 \quad k = 1, \dots, 4M + 1
\end{aligned}$$

This linear programming problem has a feasible solution:

$$\begin{aligned}
v_k &= P(T = t_k, I = 1 | X = 0) \quad k = 1, \dots, M, \\
v_{M+k} &= P(T = t_k, I = 0 | X = 0) \quad k = 1, \dots, M, \\
v_{2M+k} &= P(T = t_k, I = 1 | X = 1) \quad k = 1, \dots, M, \\
v_{3M+k} &= P(T = t_k, I = 0 | X = 1) \quad k = 1, \dots, M, \\
v_{4M+1} &= 1, \\
p(s_1, s_2) &= 0 \quad \text{for all } (s_1, s_2)
\end{aligned}$$

and the optimal function value in (11) is 0 if the equations (6)–(10) have a solution and it is strictly negative otherwise.

It is clear that the approach generalizes to the case where there are more than two latent failure times, and to the case where the (vector of) explanatory variable(s) takes more than two values.

Theorem 2 (and the accompanying corollary) of the appendix establishes that $\hat{f}(a, b)$ converges uniformly to $f(a, b)$ where the former has been defined by the same linear programming problem but with all the probabilities, P , replaced by consistent estimates. Moreover, the uniform rate of convergence equals that of \hat{P} to P . It therefore follows by the argument in Manski and Tamer (2002) that the identified region can be consistently estimated by the set of parameter values, (a, b) , such that $\hat{f}(a, b) \geq \max \hat{f} - \varepsilon_n$ where ε_n is some sequence that converges to 0 more slowly than the rate of convergence of \hat{P} .

The consistency argument with the corresponding rate of convergence is quite generic. For the particular problem studied in this paper, it is possible to establish additional results based on the following Lemma which is proved in the Appendix.

Lemma 1 *The functions $f(a, b)$ and $\hat{f}(a, b)$ are both piecewise constant over the same finite number of regions.*

Lemma 1 essentially makes the parameter space discrete with a finite number of elements.

Note that the setup in (11) forces one to underestimate all the probabilities. While this does not affect the consistency of the resulting estimator of α and β , it may be intuitively unappealing. It might therefore be more attractive to consider the linear programming problem

$$f(a, b) = \max_{\{v_i\}, \{u_i\}, \{p(\cdot, \cdot)\}} \sum - (v_i + u_i) \quad (12)$$

subject to

$$\begin{aligned} v_k - u_k + \sum_{\substack{t_k < s_1 < t_{k+1} \\ s_2 > s_1}} p(s_1, s_2) &= P(T = t_k, I = 1 | X = 0) \quad k = 1, \dots, M, \\ v_{M+k} - u_{M+k} + \sum_{\substack{t_k < s_2 < t_{k+1} \\ s_1 > s_2}} p(s_1, s_2) &= P(T = t_k, I = 0 | X = 0) \quad k = 1, \dots, M, \\ v_{2M+k} - u_{2M+k} + \sum_{\substack{t_k < as_1 < t_{k+1} \\ bs_2 > as_1}} p(s_1, s_2) &= P(T = t_k, I = 1 | X = 1) \quad k = 1, \dots, M, \\ v_{3M+k} - u_{3M+k} + \sum_{\substack{t_k < bs_2 < t_{k+1} \\ as_1 > bs_2}} p(s_1, s_2) &= P(T = t_k, I = 0 | X = 1) \quad k = 1, \dots, M, \\ v_{4M+1} - u_{4M+1} + \sum_{s_1, s_2} p(s_1, s_2) &= 1, \\ p(s_1, s_2) &\geq 0 \quad \text{for all } (s_1, s_2), \\ u_i, v_i &\geq 0 \quad k = 1, \dots, 4M + 1. \end{aligned}$$

The disadvantage of this approach is that it increases the dimensionality of the linear programming problem. In the application below, we will therefore focus on the first formulation.

5 Extensions

5.1 No assumption is made on the effect of X on T_2 .

It is relatively straightforward to establish bounds for a in the case where one makes no assumption on the effect of X on T_2 . Specifically, suppose that

$$(T^*, I) = \begin{cases} (\min\{S_1, S_2\}, 1\{S_1 < S_2\}) & \text{for } X = 0, \\ (\min\{\alpha S_1, \tilde{S}_2\}, 1\{\alpha S_1 < \tilde{S}_2\}) & \text{for } X = 1, \end{cases}$$

where (S_1, S_2, \tilde{S}_2) is independent of X . The identified region for α is the set of a 's such that there exist $p(s_1, s_2)$ and $\tilde{p}(s_1, s_2)$ satisfying

$$\begin{aligned}
\sum_{\substack{t_k < s_1 < t_k+1 \\ s_2 > s_1}} p(s_1, s_2) &= P(T = t_k, I = 1 | X = 0), \\
\sum_{\substack{t_k < s_2 < t_k+1 \\ s_1 > s_2}} p(s_1, s_2) &= P(T = t, I = 0 | X = 0) \\
\sum_{\substack{t_k < a s_1 < t_k+1 \\ s_2 > a s_1}} \tilde{p}(s_1, s_2) &= P(T = t_k, I = 1 | X = 1), \\
\sum_{\substack{t_k < s_2 < t_k+1 \\ a s_1 > s_2}} \tilde{p}(s_1, s_2) &= P(T = t, I = 0 | X = 1) \\
\sum_{s_1, s_2} p(s_1, s_2) &= 1, \quad \sum_{s_1, s_2} \tilde{p}(s_1, s_2) = 1, \\
\sum_{s_2} p(s_1, s_2) &= \sum_{s_2} \tilde{p}(s_1, s_2) \\
p(s_1, s_2) &\geq 0, \quad \tilde{p}(s_1, s_2) \geq 0
\end{aligned}$$

where the last set of equality constraints captures the constraint that the marginal distribution of S_1 should be the same whether it is calculated from the distribution of (S_1, S_2) or from the distribution of (S_1, \tilde{S}_2) . These equations again have the structure of the constraints of a linear programming problem.

As in section 4.1, one can estimate the identified region as a set of maximizers of a function that is defined as the optimal function value for a linear programming problem.

5.2 Counterfactuals

The explanatory variable, X , is often a time-dummy. In that case, it is natural to ask what the distribution of T would have been if only the distribution of T_1 had changed.

Consider for example the setup on section 4.1 and define

$$\tilde{T}_1^* = \min \{\alpha S_1, S_2\}$$

This is the duration that one would observe if X has the hypothesized effect on the first latent duration but has no effect on the second duration. This could then be compared to the distribution of T^* given $X = 1$ in order to find the effect that X has on T through T_2 alone (keeping the distribution of T_1 where it would be when $X = 1$).⁸ It might also be interesting to know the effect

⁸Other effects of this type could be considered. For example, one could compare the distribution of

$$\min \{S_1, \beta S_2\}$$

of eliminating a risk altogether. This suggests comparing the distribution of T^* given $X = 1$ to the distribution of

$$\tilde{T}_2^* = \alpha S_1$$

(or the distribution of T^* given $X = 0$ to the distribution of S_1).

Unfortunately, such an exercise is not literally possible if T^* is grouped. In that case one can only get the distribution of the grouped version of T^* given $X = 0$ or given $X = 1$. It is therefore natural to also consider the distribution of the grouped version of \tilde{T}_1^* or T_2^* . This is the equivalent of considering the distribution function for \tilde{T}^* at the points t_1, t_2, \dots etc.

For a given α and β and a given $p(\cdot, \cdot)$ we have

$$\begin{aligned} P(\tilde{T}_1^* < t_k) &= P(\min\{\alpha S_1, S_2\} < t_k) \\ &= \sum_{s_1 < t_k/\alpha \text{ or } s_2 < t_k} p(s_1, s_2) \end{aligned}$$

The last expressions are affected by the fact that the points of support are not uniquely determined. Before proceeding, it is therefore necessary to consider every polygon depicted in solid lines in the third graph of Figure 1 and determine whether the location of a point within the region changes whether the event $\{s_1 < t_k/\alpha \text{ or } s_2 < t_k\}$ occurs. If it does⁹, then one must place two points of support in the region corresponding to whether or not $\{s_1 < t_k/\alpha \text{ or } s_2 < t_k\}$.

One can then calculate population bounds on $P(\tilde{T}_1^* < t_k)$ by minimizing and maximizing (over a and b) the function $\sum_{s_1 < t_k/\alpha, s_2 < t_k} p(s_1, s_2)$ subject to (6)–(10). Unfortunately, the sample analog of this (which replaces $P(T = t_k, I = i | X = x)$ by $\hat{P}(T = t_k, I = i | X = x)$) will not produce a consistent estimator of the upper and lower bounds on $P(\tilde{T}^* < t_k)$. The reason is that there is no guarantee that the sample version of (6)–(10) will have a solution for any value of a or b .

It is also not possible to estimate the upper and lower bounds by referring to the solution to (11). The reason for this is that for a given (a, b) , the solution for $p(\cdot, \cdot)$ need not be unique. However, this suggests constructing consistent estimators for the upper and lower bounds as follows. Let $\hat{\Theta}$ be the set of maximizers of

$$f(a, b) = \max_{\{v_i\}, \{p(\cdot, \cdot)\}} \sum -v_i \quad (13)$$

to the distribution of T^* given $X = 0$ in order to find the effect that X has on T through T_2 alone (keeping the distribution of T_1 where it would be when $X = 0$).

⁹For the polygons between the two lines, $s_1 = s_2$ and $as_1 = bs_2$, the statement $\{s_1 < t_k/\alpha \text{ or } s_2 < t_k\}$ is never ambiguous as one considers different points in the polygon. Likewise, for polygons located above both of the lines and located entirely above $s_2 = t_k$, it does not matter which point one considers in the polygon as all of the points will lead to $\{s_1 > t_k/\alpha \text{ and } s_2 > t_k\}$. The same is true for polygons located to the right of both of the lines and located entirely to the right of $s_1 = t_k/\alpha$. This means that the ambiguity affects only a relatively small number of polygons.

subject to

$$v_k + \sum_{\substack{t_k < s_1 < t_k + 1 \\ s_2 > s_1}} p(s_1, s_2) = \hat{P}(T = t_k, I = 1 | X = 0) \quad k = 1, \dots, M, \quad (14)$$

$$v_{M+k} + \sum_{\substack{t_k < s_2 < t_k + 1 \\ s_1 > s_2}} p(s_1, s_2) = \hat{P}(T = t_k, I = 0 | X = 0) \quad k = 1, \dots, M, \quad (15)$$

$$v_{2M+k} + \sum_{\substack{t_k < as_1 < t_k + 1 \\ bs_2 > as_1}} p(s_1, s_2) = \hat{P}(T = t_k, I = 1 | X = 1) \quad k = 1, \dots, M, \quad (16)$$

$$v_{3M+k} + \sum_{\substack{t_k < bs_2 < t_k + 1 \\ as_1 > bs_2}} p(s_1, s_2) = \hat{P}(T = t_k, I = 0 | X = 1) \quad k = 1, \dots, M, \quad (17)$$

$$v_{4M+1} + \sum_{s_1, s_2} p(s_1, s_2) = 1, \quad (18)$$

$$p(s_1, s_2) \geq 0 \quad \text{for all } (s_1, s_2), \quad (19)$$

$$v_i \geq 0 \quad k = 1, \dots, 4M + 1 \quad (20)$$

and let \hat{f} be the optimal function value. The consistent estimators of the upper bound on $P(\tilde{T}^* < t_k)$ is then obtained by maximizing $g(a, b)$ over (a, b) in $\hat{\Theta}$ where

$$g(a, b) = \max_{\{v_i\}, \{p(\cdot, \cdot)\}} \sum_{s_1 < t_k/a, s_2 < t_k} p(s_1, s_2)$$

subject to (14)–(20) and

$$\sum -v_i = \hat{f}$$

The consistent estimators of the lower bound on $P(\tilde{T}^* < t_k)$ is obtained by minimizing $g(a, b)$ over (a, b) in $\hat{\Theta}$ where

$$g(a, b) = \min_{\{v_i\}, \{p(\cdot, \cdot)\}} \sum_{s_1 < t_k/a, s_2 < t_k} p(s_1, s_2)$$

subject to the same constraints.

The same setup can be used to construct estimates of the upper and lower bounds for objects related to life expectancy. Specifically consider the censored duration until death, $\min\{T, M\}$, rounded down to the nearest integer. The upper bound for this is obtained by maximizing $g(a, b)$ over (a, b) in $\hat{\Theta}$ where

$$g(a, b) = \max_{\{v_i\}, \{p(\cdot, \cdot)\}} \sum \text{int}\{\min a \cdot s_1, s_2, M\} \cdot p(s_1, s_2)$$

subject to the same constraints. The lower bound is obtained by minimizing the same function subject to the same constraints.

5.3 Exclusion Restrictions

Exclusion restrictions are sometimes useful in improving identification. One way to model an exclusion restriction in the competing risks model is to assume that the explanatory variable X is independent of one of the latent durations

$$(T^*, I) = \begin{cases} (\min \{S_1, S_2\}, 1 \{S_1 < S_2\}) & \text{for } X = 0, \\ (\min \{\tilde{S}_1, S_2\}, 1 \{\tilde{S}_1 < S_2\}) & \text{for } X = 1, \end{cases}$$

This model generalizes the competing risks model considered by, for example, Faraggi and Korn (1996), and it is in the spirit of many econometric models in which exclusion restrictions are used to obtain point-identification.

In this section, we will discuss how to obtain bounds on difference in the distribution functions for S_1 and \tilde{S}_1 . This is essentially done as in the same way that the Peterson bounds were constructed, but with the added restriction that the marginal distribution for S_2 is the same in the two subsamples given by $X = 0$ and $X = 1$.

Suppose that we are interested in bounding $P(T_1 \leq t | X = 1) - P(T_1 \leq t | X = 0)$ for some t . In this case, the relevant points of support are given in Figure 2.¹⁰ Most of the points of support are self-explanatory. There are, however, two main differences relative to the points of support in the first panel of Figure 1. The first is that for each region that includes $T_1 = t$ in its interior, one must allow for one point to the left of t and one to the right. The second complication is that for each region, one must allow for a point of support corresponding to each of the discrete values of T_2 that fall in the region. This is needed because one needs these to enforce the restriction that the marginal distributions of T_2 are the same in the two periods. Except for that, the points of support are as they would be in the first panel of Figure 1.

The lower bound for $P(T_1 \leq t | X = 1) - P(T_1 \leq t | X = 0)$ is then the value of

$$\min_{\{v_i\}, \{p(\cdot, \cdot)\}, \{\tilde{p}(\cdot, \cdot)\}} \sum_{s_1 \leq t} \tilde{p}(s_1, s_2) - \sum_{s_1 \leq t} p(s_1, s_2)$$

subject to

$$\begin{aligned} \sum_{\substack{t_k < s_1 < t_k + 1 \\ s_2 > s_1}} p(s_1, s_2) &= P(T = t_k, I = 1 | X = 0), & \sum_{\substack{t_k < s_2 < t_k + 1 \\ s_1 > s_2}} p(s_1, s_2) &= P(T = t, I = 0 | X = 0) \\ \sum_{\substack{t_k < s_1 < t_k + 1 \\ s_2 > s_1}} \tilde{p}(s_1, s_2) &= P(T = t_k, I = 1 | X = 1), & \sum_{\substack{t_k < s_2 < t_k + 1 \\ s_1 > s_2}} \tilde{p}(s_1, s_2) &= P(T = t, I = 0 | X = 1) \end{aligned}$$

¹⁰Figure 2 is drawn for the case where the observations are censored after 9 periods and $t = 5.5$.

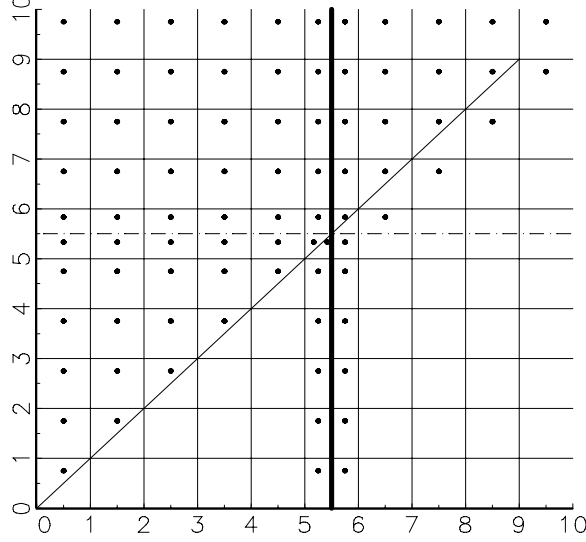


Figure 2: Illustration of Points of Support Necessary to Deal with Exclusion Restrictions

$$\sum_{s_1, s_2} p(s_1, s_2) = 1, \quad \sum_{s_1, s_2} \tilde{p}(s_1, s_2) = 1, \quad \sum_{s_1} p(s_1, s_2) = \sum_{s_1} \tilde{p}(s_1, s_2)$$

$$p(s_1, s_2) \geq 0, \quad \tilde{p}(s_1, s_2) \geq 0$$

and the upper bound for $P(T_1 \leq t | X = 1) - P(T_1 \leq t | X = 0)$ is then the value of

$$\max_{\{v_i\}, \{p(\cdot, \cdot)\}, \{\tilde{p}(\cdot, \cdot)\}} \sum_{s_1 \leq t} \tilde{p}(s_1, s_2) - \sum_{s_1 \leq t} p(s_1, s_2)$$

subject to the same constraints.

As in section 5.2, there is no guarantee that the sample analogs of these will be consistent estimators of the lower and upper bounds for $P(T_1 \leq t | X = 1) - P(T_1 \leq t | X = 0)$ as the sample analogs of the constraints may have no solution. To derive consistent estimators of these, first define \hat{f} by

$$\hat{f} = \max_{\{v_i\}, \{p(\cdot, \cdot)\}, \{\tilde{p}(\cdot, \cdot)\}} \sum -v_i$$

subject to

$$v_k + \sum_{\substack{t_k < s_1 < t_k + 1 \\ s_2 > s_1}} p(s_1, s_2) = \hat{P}(T = t_k, I = 1 | X = 0),$$

$$v_{k+M} + \sum_{\substack{t_k < s_2 < t_k + 1 \\ s_1 > s_2}} p(s_1, s_2) = \hat{P}(T = t, I = 0 | X = 0),$$

$$v_{k+2M} + \sum_{\substack{t_k < s_1 < t_k + 1 \\ s_2 > s_1}} \tilde{p}(s_1, s_2) = \hat{P}(T = t_k, I = 1 | X = 1),$$

$$v_{k+3M} + \sum_{\substack{t_k < s_2 < t_k + 1 \\ s_1 > s_2}} \tilde{p}(s_1, s_2) = \hat{P}(T = t, I = 0 | X = 1),$$

$$\begin{aligned}
v_{1+4M} + \sum_{s_1, s_2} p(s_1, s_2) &= 1, & v_{2+4M} + \sum_{s_1, s_2} \tilde{p}(s_1, s_2) &= 1, \\
\sum_{s_1} p(s_1, s_2) &= \sum_{s_1} \tilde{p}(s_1, s_2), & p(s_1, s_2) &\geq 0, & \tilde{p}(s_1, s_2) &\geq 0
\end{aligned}$$

This has a feasible solution defined, for example, by setting $p(s_1, s_2) = \tilde{p}(s_1, s_2) = 0$ for all (s_1, s_2) .

The consistent estimator of the lower bound for $P(T_1 \leq t | X = 1) - P(T_1 \leq t | X = 0)$ is then the value of

$$\min_{\{v_i\}, \{p(\cdot, \cdot)\}, \{\tilde{p}(\cdot, \cdot)\}} \sum_{s_1 \leq t} \tilde{p}(s_1, s_2) - \sum_{s_1 \leq t} p(s_1, s_2)$$

subject to the same constraints and

$$\hat{f} = \sum -v_i$$

Similarly, the upper bound for $P(T_1 \leq t | X = 1) - P(T_1 \leq t | X = 0)$ is found by maximizing $\sum_{s_1 \leq t} \tilde{p}(s_1, s_2) - \sum_{s_1 \leq t} p(s_1, s_2)$.

5.4 Bounds with Continuous Covariates or Non-grouped Durations

In the discussion above, we focused on the case where the explanatory variable X is discrete and the durations are grouped. This is the case in which the competing risk model with the parametric assumptions is most obviously not identified, and it therefore represents a worst-case scenario. On the other hand, it is also a case in which all the observed variables have a discrete distribution. This is essential for the simple approach taken above.

Each of the two complications, discrete covariates and grouped durations, violate the assumptions in for example Heckman and Honoré (1989) or Abbring and van den Berg (2003). It is therefore not known whether the parameters of the resulting models are point-identified. In this section we demonstrate that it is in principle easy to derive expressions for the identified region for those parameters (whether or not this is a single point).

First assume that X is continuous and the durations are grouped. If the model is

$$(T^*, I) = (\min\{\alpha(X) S_1, \beta(X) S_2\}, 1\{\alpha(X) S_1 < \beta(X) S_2\})$$

with the normalization $\alpha(0) = \beta(0) = 1$, then the identified region for $(\alpha(\cdot), \beta(\cdot))$ is the set of

functions $(a(\cdot), b(\cdot))$ such that there exists $p(s_1, s_2)$ satisfying

$$\int_{t_k/a(X)}^{t_{k+1}/a(X)} \int_{a(X)s_1/b(X)}^{\infty} p(s_1, s_2) ds_2 ds_1 = P(T = t_k, I = 1 | X), \quad (21)$$

$$\int_{t_k/b(X)}^{t_{k+1}/b(X)} \int_{b(X)s_2/a(X)}^{\infty} p(s_1, s_2) ds_1 ds_2 = P(T = t_k, I = 0 | X), \quad (22)$$

$$\int \int p(s_1, s_2) ds_1 ds_2 = 1, \quad (23)$$

$$p(s_1, s_2) \geq 0 \quad (24)$$

for all values of X (where the first four equations hold for all $k = 1, \dots, M$). The identified region for $(\alpha(\cdot), \beta(\cdot), p(\cdot, \cdot))$ can also be expressed as

$$\begin{aligned} \arg \min_{a(\cdot), b(\cdot), p(\cdot, \cdot)} E & \left[\sum_k \left(\int_{t_k/a(X)}^{t_{k+1}/a(X)} \int_{a(X)s_1/b(X)}^{\infty} p(s_1, s_2) ds_2 ds_1 g(X) - P(T = t_k, I = 1 | X) g(X) \right)^2 \right. \\ & \left. + \sum_k \left(\int_{t_k/b(X)}^{t_{k+1}/b(X)} \int_{b(X)s_2/a(X)}^{\infty} p(s_1, s_2) ds_1 ds_2 g(X) - P(T = t_k, I = 0 | X) g(X) \right)^2 \right] \end{aligned}$$

subject to $\int \int p(s_1, s_2) ds_1 ds_2 = 1$ and $p(s_1, s_2) \geq 0$ where $g(\cdot)$ is a positive weighting function. As discussed in Honoré and Tamer (2004), this can be turned into a feasible estimator of the identified region of (a, b) by replacing terms like $P(T = t_k, I = 1 | X)$ with the nonparametric estimates and replacing a, b and p with approximations. The weighting function $g(\cdot)$ is useful because it can be used to control for the fact that $P(T = t_k, I = 1 | X)$ will be imprecisely estimated in the tails of the distribution of X . In particular, it is straightforward to prove consistency of the estimator of the identified region for (α, β) if one uses $g(\cdot)$ to be the estimated density of X . Parametric restrictions on $\alpha(\cdot)$ and $\beta(\cdot)$ can be incorporated by minimizing the function above, subject to those restrictions.

Next consider the case where X is discrete with two points of support and the durations are not grouped. In this case, the identified region is given by the set of (a, b) for which there exists

$p(s_1, s_2)$ satisfying

$$\begin{aligned}
\int_t^\infty \int_{s_1}^\infty p(s_1, s_2) ds_2 ds_1 &= P(T > t, I = 1 | X = 0) \\
\int_t^\infty \int_t^{s_1} p(s_1, s_2) ds_2 ds_1 &= P(T > t, I = 0 | X = 0) \\
\int_{t/a}^\infty \int_{as_1/b}^\infty p(s_1, s_2) ds_2 ds_1 &= P(T > t, I = 1 | X = 1) \\
\int_{t/a}^\infty \int_{t/b}^{as_1/b} p(s_1, s_2) ds_2 ds_1 &= P(T > t, I = 0 | X = 1) \\
\int \int p(s_1, s_2) ds_1 ds_2 &= 1, \\
p(s_1, s_2) &\geq 0
\end{aligned}$$

This can also be expressed as the solutions to a population optimization problem,

$$\begin{aligned}
&\min_{a(\cdot), b(\cdot), p(\cdot, \cdot)} \int_0^\infty \left(\int_t^\infty \int_{s_1}^\infty p(s_1, s_2) ds_2 ds_1 - P(T > t, I = 1 | X = 0) \right)^2 dt \\
&+ \int_0^\infty \left(\int_t^\infty \int_t^{s_1} p(s_1, s_2) ds_2 ds_1 - P(T > t, I = 0 | X = 0) \right)^2 dt \\
&+ \int_0^\infty \left(\int_{t/a}^\infty \int_{as_1/b}^\infty p(s_1, s_2) ds_2 ds_1 - P(T > t, I = 1 | X = 1) \right)^2 dt \\
&+ \int_0^\infty \left(\int_{t/a}^\infty \int_{t/b}^{as_1/b} p(s_1, s_2) ds_2 ds_1 - P(T > t, I = 0 | X = 1) \right)^2 dt
\end{aligned}$$

subject to $\int \int p(s_1, s_2) ds_1 ds_2 = 1$ and $p(s_1, s_2) \geq 0$. This can be turned into a feasible estimator of the identified region of (a, b) by replacing terms like $P(T > t, I = 1 | X = 0)$ with the nonparametric estimates and replacing p by a sieve approximation.

6 The Change between 1970 and 2000 in the Mortality from Cancer and Cardiovascular Disease

In this section, we apply the methods described above to estimate the trends in disease-specific mortality between 1970 and 2000.

6.1 Results assuming independence

As a baseline, we start by constructing bounds under the commonly used assumption of independence. In doing so, we assume that the time dummy has a (different) multiplicative effect on the duration until death for both cancer and CVD. In order to provide a fair comparison with the rest

of our results, we use the identical estimation method, except that we estimate the bounds for the parameters separately rather than jointly. Except for the fact that we account for the grouping, this amounts to estimating a standard accelerated failure time model for each survival time. For details on the estimation, see the appendix. The conclusions from this estimation should be qualitatively similar to the conclusions we reach by looking at the raw hazard rates: the main differences here are that improvements are expressed in terms of increases in the time until death rather than in decreases in the hazard rates; that we impose a multiplicative functional form; and that we treat the data as grouped.

We compute bounds for four demographic groups separately, and for three different periods, 1970 to 1980, 1970 to 1990, and 1970 to 2000. Recall that if the duration until death has not changed since 1970, then we will find bounds around one, i.e. the duration until death in 1970 will be identical to the duration until death in a later period. Bounds above one will signal improvements. The results are in Table 2. Not surprisingly the results show large improvements from 1970 to 2000 in CVD for all groups: the duration until death from CVD increased between 30 to 40% relative to 1970. On the other hand, we find a very small, albeit positive, improvement in cancer for all groups. The survival until death from cancer increased by about 6% for white men during the same period, by about 9% for white women, and it increased by about 2% for black men and women.

For completeness, and for future reference, Table 2 also includes the results for cancer exclusive of lung cancer and from lung cancer alone. These are given in the last two rows for each panel.

6.2 Main Results

We now present our main results which construct bounds without assuming independence, as in section 4.1.¹¹ We do assume that the potential duration to death from other causes is independent of the potential duration until death from cancer and the potential duration until death from CVD.

The results are in Table 3. For all groups we find that the CVD duration increased substantially from 1970 to 2000, by about 40% for white males, 33% for blacks and 24% for white females. This increase was not concentrated in a single decade but was rather constant.

Age until death from cancer also increased for all groups during this period. This increase was about 10% for males and 15-20% for women by 2000, certainly smaller than the percentage increases

¹¹We are using estimates of the cause-specific probabilities of dying at different ages. The theory presented earlier therefore requires one to define the interval estimates as the set of parameter values for which the function value is within some ε_n of its maximum. Since measurement errors are likely to be more important than estimation uncertainty, we ignore this issue in this application.

for CVD, but not negligible. However for white males the increase was mostly concentrated in the 1990s; from 1970 to 1990 the increases were small, about 2 to 4%. The same is not true for females, who saw some significant improvements in every decade. For black males there were very small improvements in each decade.

We compare these results with those we presented in the previous section (Table 2). The coefficients for CVD are similar with or without independence, especially for white men, but the estimated improvements are larger when we do not assume independence. On the other hand the coefficients for cancer are much larger when we do not assume independence: the improvements more than double for all groups.

Overall, these bounds support the idea that there was significant progress in cancer. Importantly note that all the bounds are tightly estimated (the range of the bounds is about 0.003 and the largest range is 0.028), and they never include one. This is true whether or not we assume independence.

6.3 Policy applications: Counterfactuals

We next use the results to answer two questions. First we ask what the contribution of cancer improvements to changes in mortality would have been either in the absence of improvements for cardiovascular disease, or with the observed improvements in CVD. We estimate these counterfactuals as described in Section 5.2. Since we have censored the data at age 80 (and the model is likely to be unreliable in the tail of the distribution), we consider the effect on the probability of surviving past age 75, and on life expectancy (censored at 80). This number can be used to evaluate the progress that had been obtained thus far. Secondly we ask what the changes in mortality would be if we could eliminate cancer as a cause of death. This will give a maximum on the value of further progress in the fight on cancer.

The results are presented in Table 4. In the first row for each group we report the actual probabilities in each period. In the next row we report (bounds for) the fitted probability of surviving past age 75 in 1970. These fitted values are based on estimation using data from 1970 and 1980, 1970 and 1990, and 1970 and 2000, respectively. This explains why one should not expect the numbers in this row to be identical, although we would have considered it as evidence against the multiplicative function form if they had differed by a great deal. Similarly in the fifth row we report the fitted probability of surviving past 75 in 1980, 1990 and 2000, respectively. Comparing the fitted values in rows two and five to the actual values in the first row reveals that the fitted values are consistently very close to the actual values. This provides some evidence that our functional form assumption is not inconsistent with the data. It is also worth noting that the fitted values are always consistently below the actual values. This should be expected since our

linear programming program always underestimates the probabilities (see the end of section 4).

In the third row of Table 4 we report (bounds for) the probability of surviving past age 75 in the absence of any progress in cancer (but including progress in CVD) and in the next row we report this probability in the absence of progress in CVD (but including progress in cancer).

In the case of white males, the probability of surviving past age 75 increased by about 19.5 percentage points, from 56.1% in 1970 to 75.6% in 2000. From row 3 we see that in the absence of cancer progress this probability would have been between 66% and 73.8% in 2000. Therefore from this vantage point progress in cancer ranges from 2 to 10.6 percentage points and accounts for somewhere between 10% to 55% of the total increase in survival.

Alternatively we can look at what the probability of survival would have been in the absence of CVD progress by looking at the fourth row. In the absence of CVD progress survival rates would have been between 57.4% and 59.5% rather than 57.4%, therefore we find that for white males cancer progress accounts for about 0-11% of the total increase in survival in this period.

The difference in the estimates is driven by the choice of baseline: the first estimate uses the later year as a baseline and therefore computes counterfactuals that allow progress in CVD. The second set of estimates use 1970 as a baseline and are computed in the absence of progress in CVD. The difference tells us about the extent to which progress in cancer (in the form of reductions in cancer mortality) cannot occur without concurrent progress in CVD (and vice-versa) which depends on the extent to which the diseases are correlated (at the extreme, if the diseases were very highly correlated eliminating one would have almost no noticeable impact on mortality because everyone would die seconds later).

Similar calculations for other groups show that cancer progress accounts for 22%-100% of the total increase in survival for white women using 2000 as a base, or 0% to 25% using 1970; for black women the range is 3.7% to 40% in 2000 and 0% to 11% in 1970. Finally for black men cancer progress accounts for 5%-59.6% of improvements in 2000, but only for 0 % to 11.8% in 1970.

It is clear that the bounds on the changes discussed above are potentially too wide as they are based on a comparison of bounds for two parameters. Alternatively we can estimate bounds on the change directly. These alternative bound estimates are in Table 5. The bounds that are estimated directly are somewhat tighter, although the general pattern and the qualitative conclusions do not change.

Table 4 also presents estimates of the effect of eliminating either cancer or CVD on the probability of surviving past age 75. These can be used to estimate upper bounds on the potential benefits of additional medical innovation in one of the two causes of death, without any changes in the alternative cause of death. One conclusion to be drawn from these results is that in 1970

reductions in CVD would have much larger impacts in the overall survival, compared to cancer for males. In 2000 improvements in cancer and CVD have comparable potential benefits. For women, the results are somewhat different in that the potential gains from improvements in cancer and CVD are comparable throughout the period.

We also use our model to calculate changes in life expectancy conditional on survival to age 45 (and with censoring at 80). The results are presented in Table 6. The results in this table are very similar to those in table 4, except that we can express changes in the survival distribution in terms of additional years of life, which can be used in cost benefit analysis if we have estimates of the dollar value of an additional year of life. For example, for white males the actual increase from 1970 to 2000 is approximately 3.3 years. In the absence of progresses in cancer, progress in CVD disease would account for 0.9 to 2.6 years. One estimate of the progress in cancer is therefore the remaining 0.7 to 2.4 years of life. Given the progress that occurred in CVD, the maximum increase that could have occurred is between 3.3 and 4.3. For white females, the comparable bound is between 1.6 and 2.8 years, so the potential gains from improvements in cancer appear smaller than those for white males. This is somewhat unintuitive since cancer is a relatively larger risk for females than for males. However this is partially due to the fact that the parameter of interest is lifetime censored at 80. Since whites females live the longest (see table 1) the censoring results in lower progress estimates for women. For this reason we prefer the estimates in Table 4.

It is worth noting that the fitted values in this table are always lower than the actual values. Again, this can be explained by the fact that we underestimate all the probabilities and they therefore do not add up to one.

7 Estimation issues

7.1 Specification checks

Because lung cancer accounts for a large fraction of cancer deaths (about 50% for men and 10% for women) and it is mostly affected by smoking behavior throughout life, we may be interested in estimating trends for all cancer except lung cancer.¹² As mentioned earlier, the results assuming independence for cancer excluding lung cancer are in Table 2. Assuming independence, the estimated trends are somewhat larger if we exclude lung cancer (around 7-9% for men), especially for women. In Table 7, we present the same bounds without assuming independence. We find much

¹²Deaths from lung cancer diminished in the 1990s because of decreases in smoking that started to take place in the 1960s and that are unrelated to progress in prevention and treatment since 1973 (Andersen, Remington, Trentham-Dietz, and Reeves (2002)).

larger improvements when we exclude lung cancer for all groups. The trends are about twice as large as those that include lung cancer, about 19 and 46% for white men and women respectively, and 9 and 45% for black men and women. Again these improvements are much larger than those in Table 2 (when we assumed independence). Because lung cancer and CVD have a common risk, smoking, it may be incorrect to include lung cancer with the third cause of death which we treat as independent. We could re-estimate non-lung cancer trends by grouping all other causes of death into the “other” category, including CVD. But notice that this would not be correct either since it estimates a single trend for all other causes of death.

This suggests that it may not be appropriate to estimate trends for cancer as a whole, but rather that it would be preferable to separate cancers. As mentioned earlier, it is conceptually straightforward to extend our method to estimate trends for more than two causes of death without assuming independence. It would also be straightforward to include additional categorical covariates. However both of these extensions are computationally difficult as both the number of constraints and number of unknown parameters in the linear programming problem increase linearly in the number of causes of death and in the number of different values of the covariates. We have therefore not pursued them here (except to the extent that we estimate separate models based on gender and race).

Interestingly, excluding or including lung cancer has only a small effect in our estimates of CVD progress. The imposition of independence does not greatly affect the trends either, even though our results do suggest that cancer and CVD are dependent. Intuitively this occurs because CVD is the largest risk. One way to understand this result is to think of dependence as a form of sample selection. The potential for sample selection to generate bias depends not only on how different the excluded sample is, but also on how (relatively) large this group is. In this sense, the potential for sample selection bias is largest for the smallest risks. In practice, these results suggest that it may not be very important to consider dependence if one is interested in CVD, but it may be extremely important for all other risks, especially for smaller ones.

Another important limitation of our estimation method is that it imposes a multiplicative effect of the time dummy on both cancer and CVD durations. Alternatively we estimate bounds for cancer that impose a multiplicative effect on cancer only (as in section 5.1). These results are presented in Table 8. In all cases, relaxing the parametric assumption for CVD results in bounds that are very large, typically ranging from about 0.5 to about 2.3. Furthermore, of the 12 bounds, only one set of bounds does not contain one (white females 1970–2000). It is therefore difficult to draw any conclusions from these results. Intuitively, this is not surprising: since CVD is the largest cause of death, imposing structure on its hazard improves estimation dramatically.

Finally, we did some of the calculations using the alternative formulation of the linear programming problem given in (12). As discussed this is intuitively more appealing as it does not systematically underestimate the probabilities. The results from this are very close to those obtained from those reported here. For example, the estimated intervals for the coefficients for the change between 1970 and 2000 for white males changed from (1.389, 1.391) to (1.392, 1.400) for CVD and from (1.134, 1.153) to (1.134, 1.142) for cancer.

7.1.1 Some Data Issues

There are several data issues in calculating age-specific mortality rates using matched data from the census and the death certificate files that are potentially problematic because they may affect our trend estimates.

Age misreporting both in the census and in death certificates are an important concern. To the extent that this error is not random, it may result in biased death rates. More importantly, these biases may have changed over time.

In the census there is evidence of age heaping: individuals ages 50 and above tend to overstate their ages by “rounding up,” which results in an unusually large population for ages ending in either 5 or 0. In our data age heaping is mostly an issue for blacks. Another important issue (that cannot be fully separated from age misreporting) is that the census undercounts certain groups of the population, especially blacks, and the undercount varies with age. Furthermore, the extent of the undercount varies with the census year (Schenker (1993)). This problem is again larger for blacks than for whites.

In the death certificates, there is also error in the age at death, but this error seems to be mostly confined to blacks over the age of 65, who tend to understate their age. There is no evidence of bias in ages among whites even for those above 85 (Hill, Preston, and Rosenwaike (2000)). The overall effect of age misreporting is to downward-bias mortality for older cohorts (Preston, Elo, and Stewart (1999)).

In the absence of additional data, there is no obvious way to correct mortality rates for these problems. Overall age misreporting appears to be a very important issue mostly among blacks. These data issues suggest that our results for blacks must be taken with caution.

Another issue is whether causes of death are correctly specified in the death certificate.¹³ More importantly the issue is whether there have been significant changes from 1970 to 2000 in the

¹³For example Welch and Black (2002) report that deaths that follow surgery from cancer are not attributed to the cancer for which surgery was performed.

accuracy with which causes of death are reported. There were two changes in the International Classification of Diseases (ICD) during our period, one in 1978 (from ICD8 to ICD9) and another in 1998 (to ICD10). These changes have affected trends in mortality rates by cause, but previous research has suggested the effects of these classification changes are small for broad causes of death such as cancer and CVD (Jemal, Ward, Anderson, and Thun (2003), Klebba (1980) and Anderson, Minio, Hoyert, and Rosenberg (2001)). Furthermore, studies that have compared the causes of death reported in the death certificate with the cause of death from an autopsy, have found that the quality of death certificate reporting has not changed much since the 1960s, except perhaps for the very old (Hoel, Ron, Carter, and Mabuchi (1993)). Overall changes in the observed causes of death have not significantly changed over time for broad causes of death.

7.1.2 Additional evidence

Our findings provide support for the claim that there has been progress in cancer, measured in terms of the increases in the underlying cause-specific duration. In this section we provide evidence from other sources consistent with our findings.

We looked for any evidence that there were indeed innovations in terms of cancer treatment during the period we study, starting in the 1970s for women and mostly in the 1990s for men. We focus on improvements for the major cancer sites (excluding lung¹⁴), namely breast, prostate, colorectal and ovarian cancer. Survival from colorectal cancer, which disproportionately affects men, has improved because of a combination of earlier detection and improved treatment at earlier stages. Standard treatment for colorectal cancer changed in 1990, following a National Institutes of Health Conference recommendation, to include a combination of 5FU and leucovorin, two previously existing drugs (NIH Consensus Conference (1990)). Although treatment for prostate cancer remains controversial, clinical trials in the 1990s showed promising effects of hormonal treatment (Howe, Wingo, Thun, Ries, Rosenberg, Feigal, and Edwards (2001)).

Improvements to treat women's cancers started earlier. Mammographies started being routinely offered in the 1970s and studies in the 1970s and 1980s showed that early detection substantially improved mortality, especially for women over 50.¹⁵ Breast cancer treatment also changed in the

¹⁴The fight against lung cancer has mostly focused on reducing tobacco consumption. This effort began with the Surgeon General Report in 1964 that first publicly announced that smoking increased the risk of lung cancer, and continues today. These efforts are reflected in the trends in lung cancer many years later. To our knowledge there is no evidence of other forms of progress in lung cancer.

¹⁵A review of the evidence by the U.S. Preventive Services Task Force is available at <http://www.ahrq.gov/clinic/3rduspstf/breastcancer/brcanrr.htm#ref4>

1980s with the dissemination of adjuvant chemotherapy, including multi-agent chemotherapy and tamoxifen, and then additional changes in treatment were implemented in the early 1990s for postmenopausal women (Mariotto, Feuer, Harlan, Wun, Johnson, and Abrams (2002)). Treatment for ovarian cancer was modified in 1986 (NIH Consensus Conference (1995)) to include surgery and chemotherapy with a platinum compound (cisplatin or carboplatin) after publication of results from randomized trials which showed their effectiveness (Omura, Blessing, Ehrlich, Miller, Yordan, Creasman, and Homesley (1986)).

In spite of the fact that this evidence is consistent with our trend estimates, it is worth keeping in mind that the trends that we estimate can also reflect changes in lifestyle and demographic characteristics, some of which may reflect prevention, and some which may be completely unrelated to scientific advances in cancer. Ultimately we cannot say with certainty that the trends we estimate are uniquely related to progress in treatment or whether they also reflect prevention and cohort effects.

8 Conclusions

In this paper we show that relatively weak parametric assumptions can dramatically improve identification in competing risks models. Using a semi-parametric framework we estimate trends for cancer mortality without assuming that other risks are independent. We make no parametric assumptions on the nature of the dependence between risks, and consider an accelerated failure time model with categorical covariates and grouped durations. Because this model is not point-identified, we estimate bounds for the effects of the categorical covariates.

We use our method to estimate changes in cancer and cardiovascular mortality since 1970. The estimated bounds for the effect of time on the duration until death for either cause are extremely tight, much tighter than the bounds one can obtain without making any assumptions at all (Peterson (1976)). Such bounds can therefore be obtained under many more situations and making fewer assumptions than the previous literature has suggested.

Previous research has estimated trends in cancer mortality by assuming independence and has found little or no progress. We find that trends in cancer show much larger improvements than previously estimated. We find that time until death from cancer increased by about 10% for white males and 20% for white women from 1970 to 2000 for all cancers, and by about 19% for white males and 40% for white women if we exclude lung cancer. These estimates are more than twice as large as estimates derived under independence. These improvements are not all due to changes in smoking for younger cohorts. Also we find that not all improvements took place in the 1990s; for

women, we find significant improvements going back to the 1970s. Although less robust, we find similar results for blacks.

References

- ABBRING, J. H., AND G. J. VAN DEN BERG (2003): “The identifiability of the mixed proportional hazards competing risks model,” *Journal of the Royal Statistical Society B*, 65, 701–710.
- AHN, H., AND J. L. POWELL (1993): “Semiparametric Estimation of Censored Selection Models with a Nonparametric Selection Mechanism,” *Journal of Econometrics*, 58(1-2), 3–29.
- AMEMIYA, T. (1985): *Advanced Econometrics*. Harvard University Press.
- ANDERSEN, L. D., P. REMINGTON, A. TRENTHAM-DIETZ, AND M. REEVES (2002): “Assessing a Decade of Progress in Cancer Control,” *The Oncologist*, 7, 200–204.
- ANDERSON, R. N., A. M. MINIO, D. L. HOYERT, AND H. M. ROSENBERG (2001): “Comparability of Causes of Death Between ICD-9 and ICD-10: Preliminary Estimates,” *National Vital Statistics Reports*, 49(2).
- BAILAR, J. C., AND H. L. GORNIK (1997): “Cancer Undefeated,” *The New England Journal of Medicine*, 336.
- BAILAR, J. C., AND B. M. SMITH (1986): “Progress against cancer?,” *The New England Journal of Medicine*, 314, 1226–1232.
- BERRINGTON, A., AND I. DIAMOND (2000): “Marriage or Cohabitation: A competing risks analysis of the first-partnership formation among the 1958 birth cohort,” *Journal of the Royal Statistical Society, Series A*, 163, 127–152.
- BOOTH, A. L., AND S. E. SATCHELL (1995): “The Hazards of doing a PhD: An Analysis of Completion and Withdrawal Rates of British PhD Students in the 1980s,” *Journal of the Royal Statistical Society, Series A*, 158(2), 297–318.
- CHIANG, C. L. (1991): “Competing Risks in Mortality Analysis,” *Annual Reviews Public Health*, 12, 281–307.
- COX, D. R. (1962): *Renewal Theory*, Muthuen’s Monographs on Applied Probability and Statistics. Methuen and Co., London.

- (1972): “Regression models and life-tables,” *J. Royal Statistical Soc. Series B-statistical Methodology*, 34.
- CROWDER, M. (2001): *Classical Competing Risks*. Chapman and Hall/CRC.
- DENG, Y., J. M. QUIGLEY, AND R. VAN ORDER (2000): “Mortgage Terminations, Heterogeneity and the Exercise of Mortgage Options,” *Econometrica*, 68(2), 275–307.
- DOLL, R. (1991): “Progress against cancer: An epidemiologic assessment,” *American Journal of Epidemiology*, 134(7).
- FARAGGI, D., AND E. L. KORN (1996): “Competing Risks with Frailty Models When Treatment Affects Only One Failure Type,” *Biometrika*, 83(2), 467–471.
- FLINN, C. J., AND J. J. HECKMAN (1982): “New Methods for Analyzing Structural Models of Labor Force Dynamics,” *Journal of Econometrics*, 18, 115–168.
- HECKMAN, J. J., AND B. E. HONORÉ (1989): “The Identifiability of the Competing Risks Model,” *Biometrika*, 76, 325–330.
- (1990): “The Empirical Content of The Roy Model,” *Econometrica*, 58, 1121–1149.
- HECKMAN, J. J., J. SMITH, AND N. CLEMENTS (1997): “Making the Most Out of Programme Evaluations and Social Experiments: Accounting for Heterogeneity in Programme Impacts,” *Review of Economic Studies*, 64, 487–535.
- HILL, M. E., S. H. PRESTON, AND I. ROSENWAIKE (2000): “Age Reporting Among White Americans Aged 85+: Results Of A Record Linkage Study,” *Demography*, 37(2).
- HOEL, D. G., E. RON, R. CARTER, AND K. MABUCHI (1993): “Influence of Death Certificate Errors on Cancer Mortality Trends,” *Journal of the National Cancer Institute*, 85(13), 1063–8.
- HONORÉ, B. E., AND E. T. TAMER (2004): “Bounds on Parameters in Dynamic Discrete Choice Models,” Princeton University.
- HOWE, H., P. WINGO, M. J. THUN, L. A. RIES, H. M. ROSENBERG, E. G. FEIGAL, AND B. K. EDWARDS (2001): “Annual Report to the Nation on the Status of Cancer (1973 through 1998), Featuring Cancers with Recent Increasing Trends,” *Journal of the National Cancer Institute*, 93(11).

- JEMAL, A., E. WARD, R. N. ANDERSON, AND M. J. THUN (2003): “Influence of Rules From the Tenth Revision of the International Classification of Diseases on U.S. Cancer Mortality Trends,” *Journal of the National Cancer Institute*, 95, 1727–1728.
- KALBFLEISCH, J. D., AND R. L. PRENTICE (1980): *The Statistical Analysis of Failure Time Data*. Wiley, New York.
- KATZ, L. F., AND B. D. MEYER (1990): “Unemployment Insurance, Recall Expectations, and Unemployment Outcomes,” *Quarterly Journal of Economics*, 105, 973–1002.
- KLEBBA, A. J. (1980): “Estimates of Selected Comparability Ratios Based on Dual Coding of the 1976 Death Certificates by the Eighth and Ninth Revision of the International Classification of Diseases,” *Monthly Vital Statistics Report*, 28(11).
- LLORCA, J., AND M. DELGADO-RODRIGUEZ (2001): “Competing Risks Analysis using Markov Chains: Impact of, Cerebrovascular and Ischaemic Heart Disease in Cancer Mortality,” *International Journal of Epidemiology*, 30, 99–101.
- LYNCH, H. T., AND A. DE LA CHAPELLE (2003): “Hereditary colorectal cancer,” *New England Journal of Medicine*, 348, 919–932.
- MANSKI, C. (1990): “Nonparametric Bounds on Treatment Effects,” *American Economic Review*, 80, 319–323.
- MANSKI, C. F. (2003): *Partial Identification of Probability Distributions*. Springer Series in Statistics.
- MANSKI, C. F., AND E. TAMER (2002): “Inference on Regressions with Interval Data on a Regressor or Outcome,” *Econometrica*, 70, 519–546.
- MARIOTTO, A., E. J. FEUER, L. C. HARLAN, L.-M. WUN, K. A. JOHNSON, AND J. ABRAMS (2002): “Trends in Use of Adjuvant Multi-Agent Chemotherapy and Tamoxifen for Breast Cancer in the United States 1975-1999,” *Journal of the National Cancer Institute*, 94(21).
- NIH CONSENSUS CONFERENCE (1990): “Adjuvant Therapy for Patients with Colon and Rectal Cancer,” *Journal of the American Medical Association*, 264, 1444–1450.
- (1995): “Ovarian Cancer, Screening, Treatment and Followup,” *Journal of the American Medical Association*, 273, 491–497.

- MEYER, B. D. (1990): “Unemployment Insurance and Unemployment Spells,” *Econometrica*, 58, 757–782.
- MOLINARI, F. (2004): “Partial Identification of Probability Distributions with Misclassified Data,” Cornell University.
- NABEL, E. G. (2003): “Genomic Medicine: Cardiovascular Disease,” *New England Journal of Medicine*, 349, 60–72.
- OMURA, G., J. A. BLESSING, C. E. EHRLICH, A. MILLER, E. YORDAN, W. T. CREASMAN, AND H. D. HOMESLEY (1986): “A Randomized Trial of Cyclophosphamide and Doxorubicin with or Without Cisplatin in Advanced Ovarian Carcinoma: A Gynecologic Oncology Group Study,” *Cancer*, 57, 1725–1730.
- PETERSON, A. V. (1976): “Bounds for a joint distribution with fixed sub-distribution functions: Application to competing risks,” *Proceedings of the National Academy of Science*, 73, 11–13.
- PRENTICE, R. L., AND L. A. GLOECKLER (1978): “Regression Analysis of Grouped Survival Data with Application to Breast Cancer Data,” *Biometrics*, 34, 57–67.
- PRESTON, S. H., I. T. ELO, AND Q. STEWART (1999): “Effects of Age Misreporting on Mortality Estimates at Older Ages,” *Population Studies*, 53(2), 165–177.
- ROTHENBERG, R. B. (1994): “Competing Mortality and Progress against Cancer,” *Epidemiology*, 5, 197–203.
- ROY, A. D. (1951): “Some Thoughts on the Distribution of Earnings,” *Oxford Economic Papers (New Series)*, 3, 135–146.
- SCHENKER, N. (1993): “Undercount in the 1990 Census: Special Section,” *Journal of the American Statistical Association*, 88(423).
- SCHOENBORN, C., P. F. ADAMS, P. M. BARNES, J. L. VICKERIE, AND J. S. SCHILLER (2004): “Health Behaviors of Adults: United States, 1999–2001,” *National Center for Health Statistics. Vital Health Stat*, 10(219).
- SEER (2004): “Surveillance, Epidemiology, and End Results (SEER) Program (www.seer.cancer.gov) SEER*Stat Database: Incidence - SEER 9 Regs Public-Use, Nov 2003 Sub (1973–2001),” *National Cancer Institute, DCCPS, Surveillance Research Program, Cancer Statistics Branch, released April 2004*.

- SLUD, E., AND D. BYAR (1988): “How Dependent Causes of Death Can Make Risk Factors Appear Protective,” *Biometrics*, 44(1), 265–269.
- TSIATIS, A. (1975): “A Nonidentifiability Aspect of the Problem of Competing Risks,” *Proceedings of the National Academy of Sciences*, 72, 20–22.
- VAUPEL, J. W., AND A. I. YASHIN (1999): “Cancer Rates over Age, Time and Place: insights from Stochastic Models of Heterogeneous Populations,” *Max Plank Institute for Demographic Research Working Paper 1999-006*.
- WELCH, H. G., AND W. C. BLACK (2002): “Are Deaths Within a Month of Cancer-Directed Surgery Attributed to Cancer?,” *Journal of the National Cancer Institute*, 94(14).
- WELCH, H. G., L. M. SCHWARTZ, AND S. WOLOSHIN (2000): “Are increasing 5-year survival rates evidence of success against cancer?,” *Journal of the American Medical Association*, 283(22), 2975–2978.
- WOHLFART, J., AND P. K. ANDERSEN (2001): “Commentary: Secular Trends in the Context of Competing, Risks,” *International Journal of Epidemiology*, 30, 102–103.
- WOOSTER, R., AND B. L. WEBER (2003): “Genomic Medicine: Breast and Ovarian Cancer,” *New England Journal of Medicine*, 348, 2339–2347.

9 Appendix

9.1 Uniform convergence of objective function for linear programming

Theorem 2 Suppose that a function $f(\theta; P)$ is defined by

$$f(\theta; P) = \max_{\{p_m\}, \{v_i\}} \sum_i -v_i \quad (25)$$

subject to

$$\begin{aligned} P(j) - \sum_{m=1}^M p_m \pi_m(j; \theta) &= v_j & \text{for all } j \in \mathcal{J} \\ 1 - \sum_{m=1}^M p_m &= v_0 \\ p_m &\geq 0 & \text{for all } 1 \leq m \leq M \\ v_j &\geq 0 & \text{for all } j \in \mathcal{J} \end{aligned}$$

where $P()$ is a vector of probabilities and the functions $\pi_m(j; \theta)$ non-negative and bounded from above. Then

$$|f(\theta; P_1) - f(\theta; P_2)| \leq \text{const} \cdot (1 - c_1) + \text{const} \cdot (1 - c_2)$$

where $c_1 = \min_j \frac{P_2(j)}{P_1(j)}$, $c_2 = \min_j \frac{P_1(j)}{P_2(j)}$ and the constants are independent of P_1 , P_2 and θ .

Proof: Consider a particular value of θ .

For that θ , the linear problem clearly has a feasible solution for any vector of probabilities, P (namely $v_j = P(j)$ for $j \in \mathcal{J}$, $v_0 = 1$ and $p_m = 0$ or $m = 1, \dots, M$). Now consider a vector of probabilities, P_1 , and let (p^1, ν^1) be the maximizer for (25) with $P = P_1$.

We will now construct a feasible solution to (25) with $P = P_2$.

Define $p_m^* \equiv c_1 z_m^1 \geq 0$ where $c_1 = \min_j \frac{P_2(j)}{P_1(j)} \leq 1$. Then $p_m^1 - p_m^* = (1 - c_1) p_m^1$

Also define

$$\begin{aligned} v_j^* &\equiv P_2(j) - \sum_{m=1}^M p_m^* \pi_m(j; \theta) \quad \text{for all } j \in \mathcal{J} \\ &= P_2(j) - c_1 \left(\sum_{m=1}^M p_m^1 \pi_m(j; \theta) \right) \\ &\geq \frac{P_2(j)}{P_1(j)} P_1(j) - \frac{P_2(j)}{P_1(j)} \left(\sum_{m=1}^M p_m^1 \pi_m(j; \theta) \right) \\ &= \frac{P_2(j)}{P_1(j)} v_j \geq 0 \end{aligned}$$

$$v_0^* \equiv 1 - \sum_{m=1}^M p_m^* = 1 - c_1 \sum_{m=1}^M p_m^1 \geq 1 - \sum_{m=1}^M p_m^1 \geq 0$$

$(\{p^*\}, \{\nu^*\})$ is a feasible solution to (25) with $P = P_2$. Moreover

$$|v_j^* - v_j| \leq |P_2(j) - P_1(j)| + M(1 - c_1) \leq (1 + M)(1 - c_1)$$

(since all the π 's are between 0 and 1) and

$$|v_0^* - v_0| \leq M(1 - c_1)$$

so $f(\theta; P_2) \geq f(\theta; P_1) - \text{const} \cdot (1 - c_1)$. Interchanging P_1 and P_2 we have $f(\theta; P_1) \geq f(\theta; P_2) - \text{const} \cdot (1 - c_2)$ with $c_2 = \min_j \frac{P_1(j)}{P_2(j)}$

We therefore have that

$$\text{const} \cdot (1 - c_2) \geq f(\theta; P_2) - f(\theta; P_1) \geq -\text{const} \cdot (1 - c_1)$$

and $|f(\theta; P_2) - f(\theta; P_1)| \leq \text{const} \cdot (1 - c_1) + \text{const} \cdot (1 - c_2)$.

Corollary 3 Consider a vector of probabilities P such that $\min_j P(j) > 0$. Then

$$\left| f(\theta; \hat{P}) - f(\theta; P) \right| = O_p(\hat{P} - P).$$

9.2 Proof of Lemma 1.

The proof is broken into two steps.

First consider a particular choice of (a, b) , (a_1, b_1) . Figure X The linear programming problem () assigns probabilities to each region of the figure for (a_1, b_1) . Now imagine that one perturbs (a, b) to (a_2, b_2) . Unless new regions are created by this or existing regions disappear, location of the regions will move continuously with (a, b) . In other words, one can (uniquely) associate the regions created by (a_2, b_2) with the regions created by the original values (a_1, b_1) (provided that no new regions were created and no existing regions eliminated). One can therefore consider the feasible solution to the linear programming problem for (a_2, b_2) that leaves the probabilities for each region at the values that were optimal for (a_1, b_1) . This will leave the objective function at (a_2, b_2) at the optimal level for (a_1, b_1) . It therefore follows that $f(a_1, b_1) \leq f(a_2, b_2)$. By interchanging (a_1, b_1) and (a_2, b_2) , it follows that $f(a_1, b_1) = f(a_2, b_2)$ for two sets of parameter values (a_1, b_1) and (a_2, b_2) as long as the change from (a_1, b_1) to (a_2, b_2) moves the regions without eliminating or creating any.

In the second step of the proof, we will now argue that as one varies (a, b) regions will be created or eliminated only a finite number of times. This will establish that f takes on only a finite number of values.

Consider the figure. The solid lines are the ones that are independent of (a, b) . Step one of the proof implies that the function value only changes when the configuration of (a, b) moves the new lines in a way that creates or eliminates region. It is convenient to reparameterize from (a, b) to $(\frac{a}{b}, b)$. As one varies $\frac{a}{b}$ from 0 to ∞ , the line $s_2 = \frac{a}{b}s_1$ will touch the intersection of a vertical and horizontal solid line only a finite number of times. Now consider a particular value of $\frac{a}{b}$ and consider the effect of varying b (holding $\frac{a}{b}$ fixed). Each value of b (along with the implied value of a) will allocate each of the points in $\{\frac{1}{b}, \frac{2}{b}, \dots, \frac{T_{\max}}{b}, \frac{1}{a}, \frac{2}{a}, \dots, \frac{T_{\max}}{a}\}$ into one of the sets $]0, 1[, [1, 1],]1, 2[, [2, 2],]2, 3[, [3, 3], \dots,]T_{\max} - 1, T_{\max}[, [T_{\max}, T_{\max}]$. Since the number of points in $\{\frac{1}{b}, \frac{2}{b}, \dots, \frac{T_{\max}}{b}, \frac{1}{a}, \frac{2}{a}, \dots, \frac{T_{\max}}{a}\}$ is finite and the number of sets into which they are allocated are finite, it follows that this particular $\frac{a}{b}$ will generate only a finite number of regions.

9.3 The Data

9.3.1 Population data

These data come from April 1st population counts from the Census Bureau, from the following sources:

1. 1970 population counts obtained from U.S. Bureau of the Census, Census of Population: 1970 General Population Characteristics Final Report PC(1)-B1 United States Summary.
2. 1980 Data was found at
<http://www.census.gov/population/estimates/nation/e80s/E8081RQI.txt>
3. 1990 data was found at
<http://www.census.gov/population/estimates/nation/e90s/E9090RMP.txt>
4. 2000 White population counts obtained from Census table PCT12A, Black population counts from table PCT12B and total population counts from PCT12. All three tables were found at the US Census Bureau's website: <http://factfinder.census.gov/servlet>

9.3.2 Death rate—causes of death classification

Deaths from cardiovascular diseases included ICD8 and ICD9 codes 390-458, and ICD10 codes G45, G46 and I00-I99. Deaths from cancer included ICD8 and ICD9 codes 140-239, and ICD10 codes C00 through D48. Lung cancer includes ICD8 and ICD9 codes 162, and ICD10 code C34. All other diseases were counted under the category “other causes of death”.

9.4 Details about the Calculations

The function value that defines the identified region was calculated over three grids.

The first grid was defined by the rectangle $\{0.90, 0.95, 1.00, \dots, 1.40\} \times \{0.90, 0.95, 1.00, \dots, 1.40\}$.

The second grid was defined by first calculation the set of maximizers over the original grid. Let θ_1^{\min} and θ_1^{\max} denote the minimum and maximum value of the first coordinate in that set and let θ_2^{\min} and θ_2^{\max} denote the minimum and maximum value of the second coordinate in the set. The second grid is then given by $\{\theta_1^{\min} - 0.05, \theta_1^{\min} - 0.04, \theta_1^{\min} - 0.03, \dots, \theta_1^{\max} + 0.08\} \times \{\theta_2^{\min} - 0.05, \theta_2^{\min} - 0.04, \theta_2^{\min} - 0.03, \dots, \theta_2^{\max} + 0.08\}$.

The third grid was defined in terms of the maximizers over the first two grid. Let θ_1^{\min} and θ_1^{\max} denote the minimum and maximum value of the first coordinate in that set and let θ_2^{\min} and θ_2^{\max}

denote the minimum and maximum value of the second coordinate in the set. The third grid is then given by $\{\theta_1^{\min} - 0.01, \theta_1^{\min} - 0.009, \theta_1^{\min} - 0.008, \dots, \theta_1^{\max} + 0.015\} \times \{\theta_2^{\min} - 0.01, \theta_2^{\min} - 0.009, \theta_2^{\min} - 0.008, \dots, \theta_2^{\max} + 0.015\}$.

The estimated identified region is then the set of maximizers of the union of the three grids. The numbers reported in the tables are the minimum and maximum values of each coordinate.

9.5 Estimation under independence.

To estimate the parameters under independence, we first estimate the marginal distribution of T_1 and T_2 using a Kaplan–Meier estimator. We then estimate a by

$$f(a) = \max_{\{v_i\}, \{p(\cdot, \cdot)\}} \sum -v_i$$

subject to

$$\begin{aligned} v_k + \sum_{\substack{t_k < s_1 < t_k + 1 \\ s_2 > s_1}} p(s_1) &= \hat{P}(T = t_k, I = 1 | X = 0) & k = 1, \dots, M, \\ v_{M+k} + \sum_{\substack{t_k < as_1 < t_k + 1 \\ bs_2 > as_1}} p(s_1) &= \hat{P}(T = t_k, I = 1 | X = 1) & k = 1, \dots, M, \\ v_{2M+1} + \sum_{s_1, s_2} p(s_1) &= 1, \\ p(s_1) &\geq 0 & \text{for all } s_1, \\ v_i &\geq 0 & k = 1, \dots, 2M + 1 \end{aligned}$$

where the points of support are determined in a way that is similar to the way we did it without independence.

b is estimated analogously.

**TABLE 1: Summary statistics by race, gender and decade
(conditional on survival to age 45)**

	1970	1980	1990	2000
White Males				
Age at death—all causes	70.43	72.00	73.62	74.70
Age at death from cardiovascular disease	71.57	72.99	74.51	75.97
Age at death from cancer	69.12	70.40	71.75	72.67
Age at death from other causes	68.18	70.96	73.32	74.17
Fraction deaths from cardiovascular disease	0.63	0.58	0.50	0.44
Fraction deaths from cancer	0.14	0.17	0.19	0.20
White Females				
Age at death—all causes	74.65	76.89	78.80	80.20
Age at death from cardiovascular disease	77.31	79.50	81.24	82.77
Age at death from cancer	68.37	70.54	72.57	73.86
Age at death from other causes	71.76	75.38	78.86	80.14
Fraction deaths from cardiovascular disease	0.62	0.59	0.51	0.45
Fraction deaths from cancer	0.17	0.19	0.19	0.18
Black Males				
Age at death—all causes	66.09	68.09	69.40	69.23
Age at death from cardiovascular disease	67.65	69.50	70.43	70.44
Age at death from cancer	66.30	67.90	69.42	69.73
Age at death from other causes	63.10	65.85	67.76	67.54
Fraction deaths from cardiovascular disease	0.56	0.51	0.46	0.43
Fraction deaths from cancer	0.14	0.18	0.21	0.21
Black Females				
Age at death—all causes	68.21	71.42	73.64	74.74
Age at death from cardiovascular disease	70.18	73.46	75.47	76.87
Age at death from cancer	64.63	67.30	69.39	70.21
Age at death from other causes	65.50	69.86	73.35	74.34
Fraction deaths from cardiovascular disease	0.61	0.56	0.51	0.46
Fraction deaths from cancer	0.15	0.18	0.20	0.19

TABLE 2: Marginal Identified Regions Assuming Independence

	1970–1980	1970–1990	1970–2000
White Males			
Coefficient on CVD	(1.126, 1.129)	(1.239, 1.250)	(1.392, 1.400)
Coefficient on Cancer	(1.001, 1.029)	(1.001, 1.029)	(1.059, 1.060)
Coef. on Cancer (excl lung)	(1.091, 1.093)	(1.126, 1.129)	(1.075, 1.076)
Coef. on Lung Cancer	(0.910, 0.911)	(0.905, 0.909)	(0.968, 0.968)
White Females			
Coefficient on CVD	(1.091, 1.093)	(1.201, 1.206)	(1.286, 1.291)
Coefficient on Cancer	(1.001, 1.029)	(1.059, 1.060)	(1.087, 1.093)
Coef. on Cancer (excl lung)	(1.091, 1.093)	(1.236, 1.250)	(1.334, 1.346)
Coef. on Lung Cancer	(0.843, 0.852)	(0.849, 0.852)	(0.840, 0.851)
Black Males			
Coefficient on CVD	(1.126, 1.129)	(1.201, 1.206)	(1.316, 1.320)
Coefficient on Cancer	(0.972, 0.999)	(0.965, 0.965)	(1.001, 1.029)
Coef. on Cancer (excl lung)	(1.084, 1.090)	(1.091, 1.093)	(1.091, 1.093)
Coef. on Lung Cancer	(0.847, 0.848)	(0.847, 0.851)	(0.847, 0.852)
Black Females			
Coefficient on CVD	(1.160, 1.166)	(1.273, 1.280)	(1.334, 1.346)
Coefficient on Cancer	(1.001, 1.029)	(0.972, 0.999)	(1.001, 1.029)
Coef. on Cancer (excl lung)	(1.059, 1.060)	(1.126, 1.129)	(1.239, 1.250)
Coef. on Lung Cancer	(0.840, 0.846)	(0.840, 0.842)	(0.851, 0.852)

**TABLE 3: Marginal Identified Regions
without Assuming Independence**

	1970–1980	1970–1990	1970–2000
White Males			
Coefficient on CVD	(1.126, 1.129)	(1.295, 1.296)	(1.389, 1.391)
Coefficient on Cancer	(1.001, 1.029)	(1.020, 1.035)	(1.134, 1.153)
White Females			
Coefficient on CVD	(1.092, 1.093)	(1.160, 1.160)	(1.236, 1.238)
Coefficient on Cancer	(1.091, 1.092)	(1.154, 1.157)	(1.201, 1.206)
Black Males			
Coefficient on CVD	(1.126, 1.129)	(1.201, 1.206)	(1.334, 1.346)
Coefficient on Cancer	(1.030, 1.034)	(1.063, 1.066)	(1.072, 1.074)
Black Females			
Coefficient on CVD	(1.158, 1.159)	(1.231, 1.235)	(1.334, 1.346)
Coefficient on Cancer	(1.096, 1.096)	(1.167, 1.172)	(1.158, 1.159)

TABLE 4: Counterfactual Probability of Surviving Age 75

	1970–1980	1970–1990	1970–2000
White Males			
Change in actual prob.	0.561 – 0.636	0.561 – 0.707	0.561 – 0.756
No Progress (fitted in 70)	(0.567, 0.567)	(0.572, 0.573)	(0.574, 0.574)
Progress in CVD	(0.634, 0.643)	(0.699, 0.719)	(0.661, 0.738)
Progress in Cancer	(0.567, 0.572)	(0.572, 0.579)	(0.574, 0.595)
Progress in Both (fitted in end year)	(0.638, 0.643)	(0.707, 0.719)	(0.756, 0.769)
Elim. CVD — Cancer as in end year	(0.638, 0.913)	(0.707, 0.906)	(0.756, 0.919)
Elim. Cancer — CVD as in end year	(0.661, 0.732)	(0.747, 0.812)	(0.785, 0.848)
Elim. CVD — Cancer as in 70	(0.634, 0.909)	(0.699, 0.906)	(0.661, 0.888)
Elim. Cancer — CVD as in 70	(0.567, 0.656)	(0.572, 0.662)	(0.574, 0.663)
White Females			
Change in actual prob.	0.733 – 0.784	0.733 – 0.820	0.733 – 0.843
No Progress (fitted in 70)	(0.736, 0.736)	(0.738, 0.738)	(0.739, 0.740)
Progress in CVD	(0.736, 0.776)	(0.738, 0.802)	(0.739, 0.820)
Progress in Cancer	(0.736, 0.751)	(0.738, 0.760)	(0.739, 0.765)
Progress in Both (fitted in end year)	(0.786, 0.787)	(0.824, 0.826)	(0.843, 0.850)
Elim. CVD — Cancer as in end year	(0.786, 0.917)	(0.824, 0.923)	(0.843, 0.930)
Elim. Cancer — CVD as in end year	(0.786, 0.865)	(0.824, 0.900)	(0.843, 0.916)
Elim. CVD — Cancer as in 70	(0.736, 0.906)	(0.738, 0.900)	(0.739, 0.900)
Elim. Cancer — CVD as in 70	(0.736, 0.824)	(0.738, 0.827)	(0.739, 0.827)
Black Males			
Change in actual prob.	0.473 – 0.540	0.473 – 0.577	0.473 – 0.634
No Progress (fitted in 70)	(0.474, 0.474)	(0.481, 0.482)	(0.486, 0.486)
Progress in CVD	(0.513, 0.542)	(0.536, 0.579)	(0.582, 0.629)
Progress in Cancer	(0.474, 0.481)	(0.481, 0.490)	(0.486, 0.504)
Progress in Both (fitted in end year)	(0.542, 0.542)	(0.587, 0.588)	(0.635, 0.647)
Elim. CVD — Cancer as in end year	(0.542, 0.872)	(0.587, 0.871)	(0.635, 0.880)
Elim. Cancer — CVD as in end year	(0.558, 0.660)	(0.621, 0.708)	(0.678, 0.768)
Elim. CVD — Cancer as in 70	(0.513, 0.872)	(0.536, 0.862)	(0.582, 0.862)
Elim. Cancer — CVD as in 70	(0.474, 0.589)	(0.481, 0.597)	(0.486, 0.599)
Black Females			
Change in actual prob.	0.586 – 0.673	0.586 – 0.713	0.586 – 0.748
No Progress (fitted in 70)	(0.594, 0.594)	(0.598, 0.598)	(0.603, 0.603)
Progress in CVD	(0.616, 0.673)	(0.598, 0.696)	(0.668, 0.740)
Progress in Cancer	(0.594, 0.604)	(0.598, 0.615)	(0.603, 0.622)
Progress in Both (fitted in end year)	(0.679, 0.683)	(0.717, 0.725)	(0.748, 0.764)
Elim. CVD — Cancer as in end year	(0.679, 0.907)	(0.717, 0.909)	(0.748, 0.916)
Elim. Cancer — CVD as in end year	(0.690, 0.774)	(0.726, 0.807)	(0.768, 0.846)
Elim. CVD — Cancer as in 70	(0.616, 0.897)	(0.598, 0.881)	(0.668, 0.891)
Elim. Cancer — CVD as in 70	(0.594, 0.691)	(0.598, 0.694)	(0.603, 0.698)

**TABLE 5: Bound on the Change in Counterfactual
Probability of Surviving Age 75**

	1970–1980	1970–1990	1970–2000
White Males			
Increase in prob. (w. impr in CVD)	(0.000, 0.027)	(0.000, 0.033)	(0.027, 0.108)
Increase in prob. (no impr in CVD)	(0.000, 0.005)	(0.000, 0.011)	(0.000, 0.022)
White Females			
Increase in prob. (w. impr in CVD)	(0.010, 0.051)	(0.021, 0.087)	(0.029, 0.111)
Increase in prob. (no impr in CVD)	(0.000, 0.015)	(0.000, 0.021)	(0.000, 0.025)
Black Males			
Increase in prob. (w. impr in CVD)	(0.000, 0.029)	(0.009, 0.051)	(0.016, 0.065)
Increase in prob. (no impr in CVD)	(0.000, 0.006)	(0.000, 0.008)	(0.000, 0.018)
Black Females			
Increase in prob. (w. impr in CVD)	(0.009, 0.067)	(0.028, 0.127)	(0.022, 0.097)
Increase in prob. (no impr in CVD)	(0.000, 0.010)	(0.000, 0.017)	(0.000, 0.019)

**TABLE 6: Counterfactual (Censored)
Life Expectancy at age 45**

	1970–1980	1970–1990	1970–2000
White Males			
Change in actual $E[T T \geq 45]$.	72.73 – 74.09	72.73 – 75.28	72.73 – 76.05
No Progress (fitted in 70)	(72.66, 72.71)	(72.61, 72.66)	(72.63, 72.68)
Progress in CVD	(73.46, 73.93)	(74.26, 74.95)	(73.62, 75.36)
Progress in Cancer	(72.66, 72.81)	(72.61, 72.78)	(72.63, 73.07)
Progress in Both (fitted in end year)	(73.89, 73.98)	(74.92, 74.98)	(75.64, 75.76)
Elim. CVD — Cancer as in end year	(73.89, 78.33)	(74.92, 78.05)	(75.64, 78.22)
Elim. Cancer — CVD as in end year	(74.15, 75.42)	(75.49, 76.49)	(76.06, 77.01)
Elim. CVD — Cancer as in 70	(73.46, 78.27)	(74.26, 78.02)	(73.62, 77.82)
Elim. Cancer — CVD as in 70	(72.66, 74.19)	(72.61, 74.14)	(72.63, 74.16)
White Females			
Change in actual $E[T T \geq 45]$.	75.67 – 76.51	75.67 – 77.09	75.67 – 77.46
No Progress (fitted in 70)	(75.65, 75.66)	(75.63, 75.65)	(75.61, 75.63)
Progress in CVD	(75.65, 76.27)	(75.63, 76.60)	(75.68, 76.83)
Progress in Cancer	(75.65, 75.91)	(75.63, 76.02)	(75.61, 76.07)
Progress in Both (fitted in end year)	(76.43, 76.47)	(76.92, 76.97)	(77.25, 77.31)
Elim. CVD — Cancer as in end year	(76.43, 78.47)	(76.92, 78.46)	(77.25, 78.54)
Elim. Cancer — CVD as in end year	(76.43, 77.85)	(76.92, 78.21)	(77.30, 78.44)
Elim. CVD — Cancer as in 70	(75.65, 78.27)	(75.63, 78.10)	(75.68, 78.08)
Elim. Cancer — CVD as in 70	(75.65, 77.23)	(75.63, 77.21)	(75.61, 77.16)
Black Males			
Change in actual $E[T T \geq 45]$.	70.69 – 71.98	70.69 – 72.76	70.69 – 73.80
No Progress (fitted in 70)	(70.68, 70.69)	(70.61, 70.67)	(70.62, 70.66)
Progress in CVD	(71.39, 71.91)	(71.51, 72.42)	(72.19, 73.26)
Progress in Cancer	(70.68, 70.82)	(70.61, 70.90)	(70.62, 70.93)
Progress in Both (fitted in end year)	(71.93, 71.97)	(72.48, 72.63)	(73.39, 73.52)
Elim. CVD — Cancer as in end year	(71.93, 77.73)	(72.48, 77.49)	(73.39, 77.54)
Elim. Cancer — CVD as in end year	(72.22, 74.00)	(72.93, 74.62)	(74.08, 75.46)
Elim. CVD — Cancer as in 70	(71.39, 77.67)	(71.51, 77.27)	(72.19, 77.28)
Elim. Cancer — CVD as in 70	(70.68, 72.73)	(70.61, 72.70)	(70.62, 72.65)
Black Females			
Change in actual $E[T T \geq 45]$.	72.83 – 74.45	72.83 – 75.15	72.83 – 75.75
No Progress (fitted in 70)	(72.71, 72.75)	(72.77, 72.80)	(72.72, 72.77)
Progress in CVD	(73.09, 74.04)	(73.01, 74.47)	(73.45, 74.96)
Progress in Cancer	(72.71, 72.98)	(72.77, 73.17)	(72.72, 73.13)
Progress in Both (fitted in end year)	(74.16, 74.26)	(74.80, 74.90)	(75.22, 75.35)
Elim. CVD — Cancer as in end year	(74.16, 77.96)	(74.80, 77.93)	(75.22, 77.89)
Elim. Cancer — CVD as in end year	(74.29, 75.90)	(74.94, 76.41)	(75.53, 76.82)
Elim. CVD — Cancer as in 70	(73.09, 77.73)	(73.01, 77.51)	(73.45, 77.51)
Elim. Cancer — CVD as in 70	(72.71, 74.56)	(72.77, 74.61)	(72.72, 74.51)

TABLE 7: Marginal Identified Regions Excluding Lung Cancer

	1970–1980	1970–1990	1970–2000
White Males			
Coefficient on CVD	(1.126, 1.129)	(1.295, 1.296)	(1.392, 1.399)
Coef. on Cancer (excl. lung)	(1.091, 1.093)	(1.039, 1.045)	(1.236, 1.249)
White Females			
Coefficient on CVD	(1.091, 1.093)	(1.201, 1.206)	(1.267, 1.269)
Coef. on Cancer (excl. lung)	(1.126, 1.129)	(1.239, 1.249)	(1.455, 1.458)
Black Males			
Coefficient on CVD	(1.126, 1.129)	(1.202, 1.206)	(1.334, 1.346)
Coef. on Cancer (excl. lung)	(1.112, 1.115)	(1.201, 1.205)	(1.118, 1.119)
Black Females			
Coefficient on CVD	(1.154, 1.157)	(1.286, 1.296)	(1.334, 1.346)
Coef. on Cancer (excl. lung)	(1.106, 1.111)	(1.143, 1.148)	(1.308, 1.319)

**TABLE 8: Marginal Identified Regions
(only Cancer multiplicative)**

	1970–1980	1970–1990	1970–2000
White Males			
Coefficient on Cancer	(0.520, 2.186)	(0.602, 2.124)	(0.654, 2.124)
White Females			
Coefficient on Cancer	(0.802, 1.610)	(0.890, 1.646)	(1.002, 1.698)
Black Males			
Coefficient on Cancer	(0.449, 2.356)	(0.484, 2.200)	(0.550, 2.332)
Black Females			
Coefficient on Cancer	(0.556, 2.284)	(0.644, 2.230)	(0.702, 2.332)

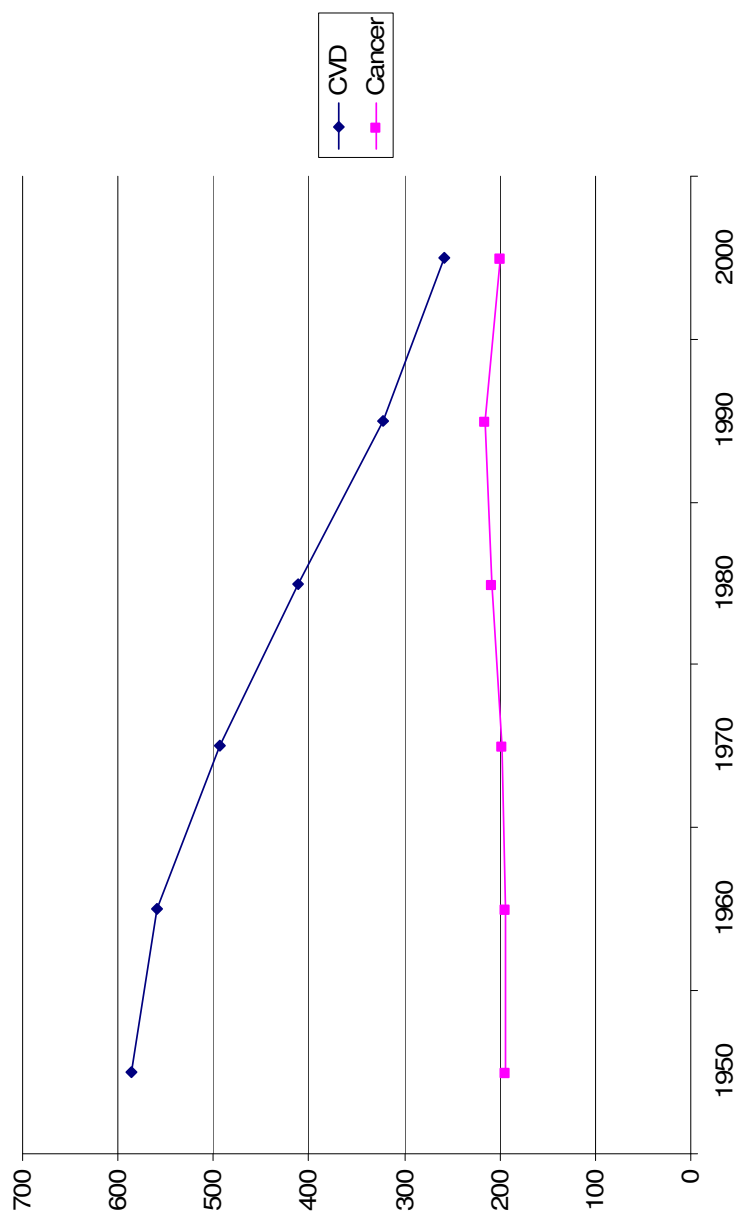


Figure 3: Trends in age-adjusted mortality 1950–2000 (all persons)

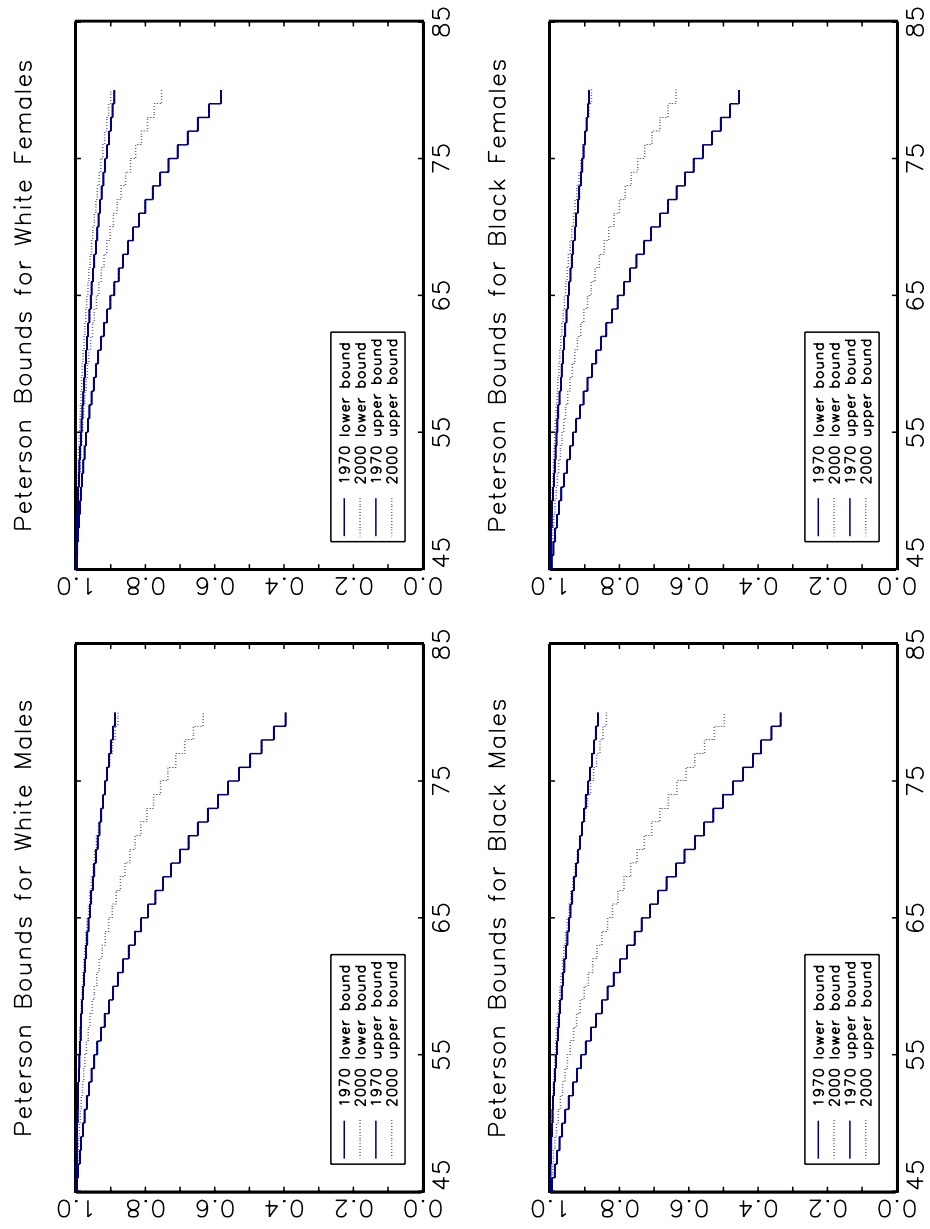


Figure 4: Peterson Bounds on the Survivor Functions in 1970 and 2000

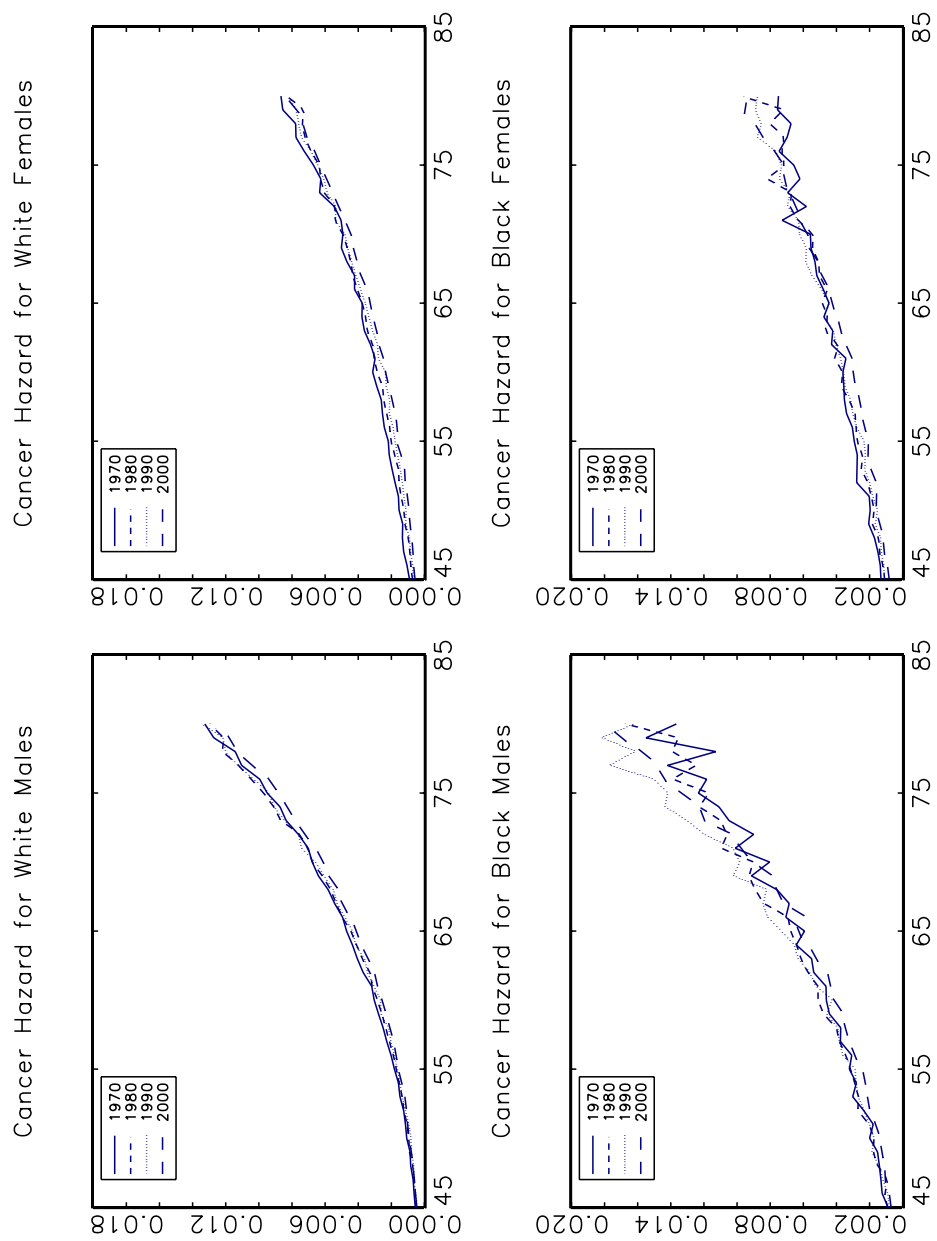


Figure 5: Hazard Rates for the Cancer

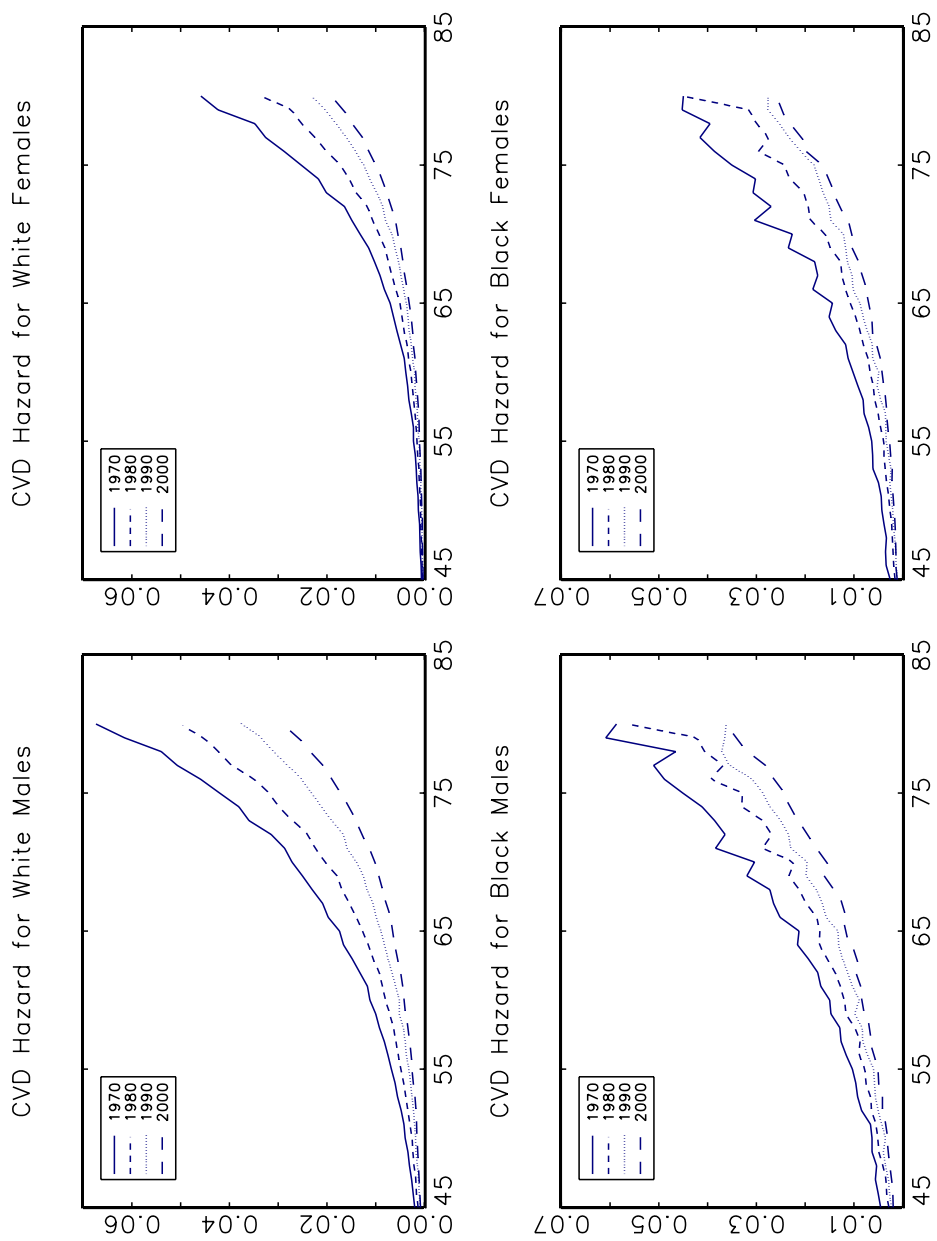


Figure 6: Hazard Rates for the CVD

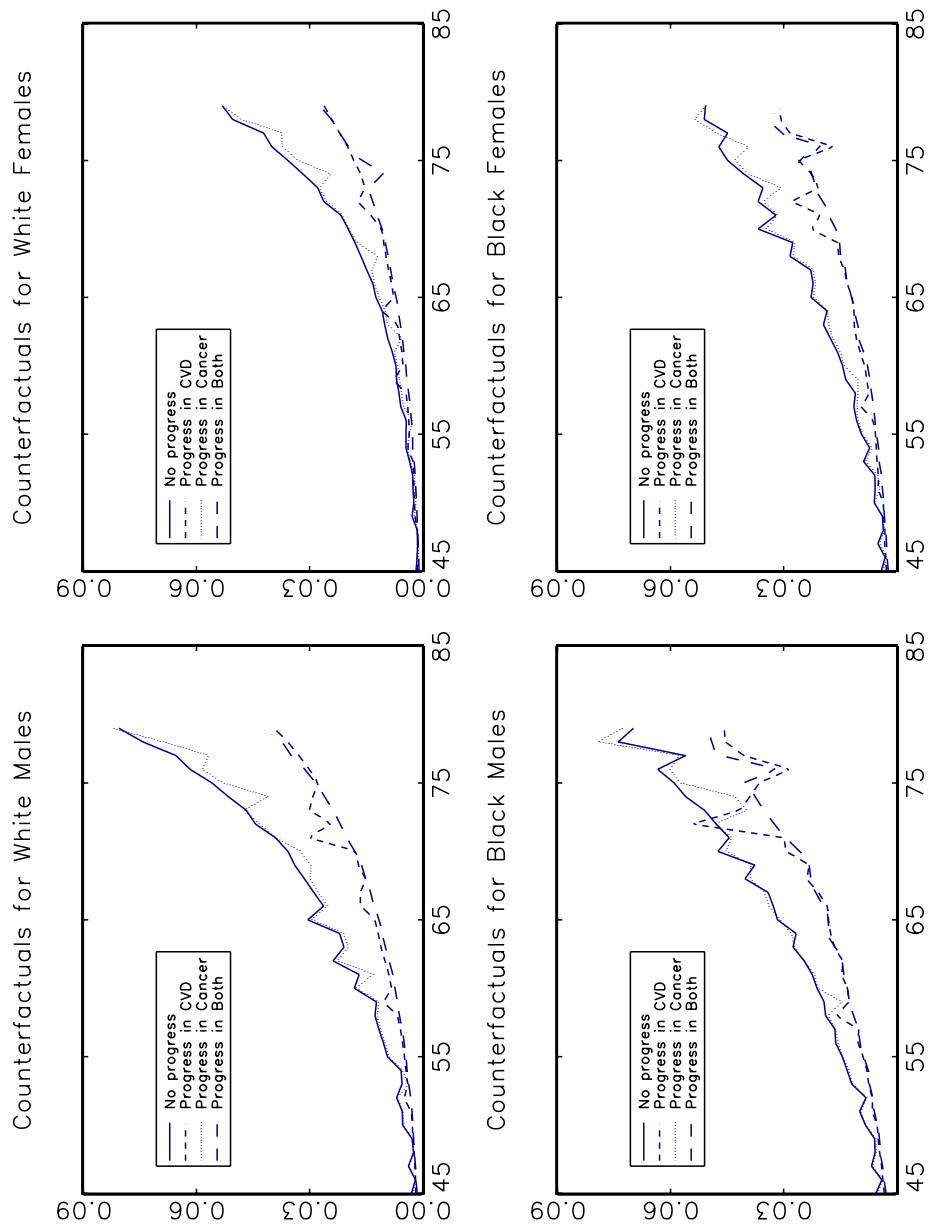


Figure 7: Counterfactual Policy Evaluations