

Forecasting United States mortality using cohort smoking histories

Haidong Wang^a and Samuel H. Preston^{b,1}

^aInstitute for Health Metrics and Evaluation, University of Washington, Seattle, 2301 5th Avenue, Seattle, WA 98121; and ^bPopulation Studies Center, University of Pennsylvania, 3718 Locust Walk, Philadelphia, PA 19104-6298

Contributed by Samuel H. Preston, November 19, 2008 (sent for review August 29, 2008)

In this paper, we introduce a recently established relationship between cohort smoking patterns and adult mortality into mortality projections for the United States. In particular, we incorporate a variable representing the intensity of smoking within a cohort into the original Lee–Carter projection model. The introduction of this variable accounts for important anomalies in the recent age/sex pattern of mortality change and enables the use of a common temporal trend of mortality change for the 2 sexes. We project age-specific mortality rates for men and women at ages 50–84 between 2004 and 2034 in the United States. Because of reductions in smoking that have already occurred or can be reliably projected, mortality is projected to decline much faster when smoking is introduced into the model.

Cigarette smoking has long been recognized as a significant factor affecting human mortality (1). The effect of smoking on mortality is clearly visible in the mortality profiles of different birth cohorts in the United States (2). These cohort patterns offer an opportunity to add important information to mortality projections. In this paper we incorporate a variable representing the prevalence of smoking among cohorts into the Lee–Carter model, the most widely used model to forecast mortality. We demonstrate the effect of this incremental information on forecasts of American mortality.

Smoking and Mortality

Although epidemiologic studies have identified cigarette smoking as an important risk factor in mortality for at least a half century, there is no consensus on how great the risk is. The large studies that have garnered the bulk of attention are based on nonrepresentative samples and have typically recorded smoking status at baseline without changing the classification of individuals as their smoking status changes. Measurement error in smoking status has the effect of reducing the estimated risk. The two largest studies, each involving more than 1 million individuals, were conducted under the auspices of the American Cancer Society. Cancer Prevention Study I, conducted between 1959 and 1972, found a ratio of mortality among current smokers relative to never-smokers ages 40–84 of 1.91 for white men and 1.46 for white women (3). Cancer Prevention Study II found ratios of 2.52 to 2.82 for men between ages 50 and 79 during 1982–1996 and ratios of 1.33–1.89 for women at these ages (4). When adjustment for smoking cessation was made for a subsample for whom this information was gathered in Cancer Prevention Study II, the range of relative risks for current smokers rose to 3.11–3.53 for men and to 2.58–3.14 for women (4). This study also produced higher estimates of the long-term effects of smoking for those who had stopped. Rogers *et al.* (5) used a nationally representative sample and found risks for current and ex-smokers that are intermediate between the 2 Cancer Prevention Studies, controlling for education, family income, and a variety of other variables.

These studies show that the risk of death from smoking is high and endures many years after smoking has ceased. In view of the high percentage of the American population that consists of current or past smokers, a percentage that reached 77% in some male cohorts (3), it is not surprising that smoking has left a distinctive

imprint on national mortality patterns. And in view of the enduring impact of smoking and its serial correlation across the life cycle, it is not surprising that much of that imprint is cohort-specific.

Smoking Patterns

The prevalence of cigarette smoking in the United States first rose and then fell in the course of the 20th century. However, the patterns were not the same for men and women (6). The sex differential in smoking prevalence also rose and fell, which was associated with widening and then narrowing sex differentials in mortality (2). Only about 6% of women smoked in 1924, a number that increased to 16% by 1929. However, during the same period, more than 50% of men smoked (6). For United States adults ages 18 and older, smoking prevalence was 56.9% for men and 28.4% for women in 1955. This sex difference of 28.5% subsequently declined to $\approx 5\%$ and has remained there since the beginning of the 1990s. In this paper we rely upon a careful reconstruction of cohort smoking information for the United States that was done by Burns *et al.* (7).

Mortality Projections

There is no universally accepted method for projecting mortality. The methods that are used can be classified into 4 categories: projections by extrapolating age-specific mortality rates; projections by reference to a model mortality scenario; projections by interpolating current mortality and a targeted mortality in the future; and projections by reference to components of mortality (8).

There have been 3 major programs of mortality forecasting in the United States: projections by the Social Security Administration (SSA), by the Census Bureau, and a series of projections using a method developed by Lee and Carter (9). All of these programs are based on observations of past trends in period levels of mortality and assumptions about whether and how those trends will be modified in the future. The Census Bureau (10) uses Lee–Carter procedures to establish long-range targets to which estimates eventually converge. The Lee–Carter method has also been recommended by advisory panels to SSA for use by the SSA (11).

The Lee–Carter model (9, 12) has many advantages: it produces an excellent fit to mortality trends; it is parsimonious in the number of parameters used; it linearizes mortality trends and thereby adds confidence to extrapolations; and it produces sensible estimates of forecast uncertainty. The basic model can be expressed as

$$\ln(M_{x,t}) = a_x + b_x k_t + \varepsilon_{x,t} \quad [1]$$

In this model, $M_{x,t}$ is the central mortality rate for age x at time t , whereas a_x describes the mortality profile by age, which is constant over time; k_t represents the temporal trend of mortality

Author contributions: H.W. and S.H.P. designed research; H.W. performed research; H.W. and S.H.P. analyzed data; and H.W. and S.H.P. wrote the paper.

The authors declare no conflict of interest.

Freely available online through the PNAS open access option.

¹To whom correspondence should be addressed. E-mail: spreston@falcon.sas.upenn.edu.

© 2009 by The National Academy of Sciences of the USA

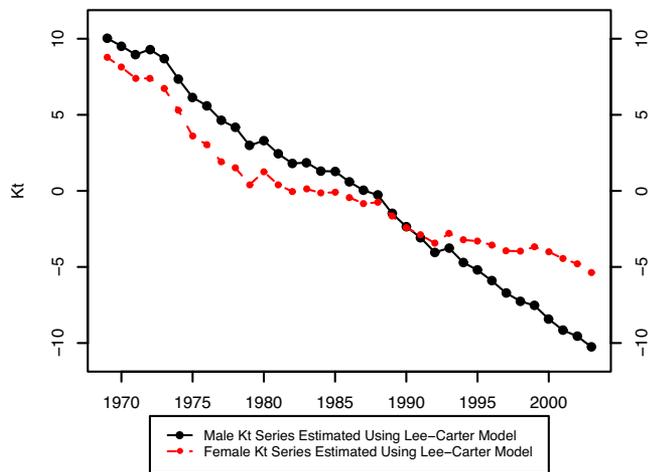


Fig. 1. K_t estimated by using the Lee-Carter Model: United States, 1969–2003.

changes over time; and b_x describes the changes in the mortality rates at age x in response to changes in k_t over time. The $\varepsilon_{x,t}$ is an error term, which depicts the age period-specific influences that are not explained by the model. In the typical application, this model is fitted to past data to obtain parameter estimates, and projections are then made by holding the a_x and b_x parameters constant and extrapolating k_t .

In their 1992 paper, Lee and Carter applied their model to the sex-combined United States population instead of applying it separately to male and female populations; they were concerned that extrapolating sex differentials in mortality would produce implausible differentials in the future. To illustrate the difficulty, we have estimated the Lee-Carter model separately for male and female United States populations between 1969 and 2003 at ages 40–84. Fig. 1 shows that the temporal trends in mortality are markedly different for males and females. In the past 35 years, male mortality has declined much more rapidly than female mortality, a pattern that may produce implausible results if projected far into the future.

Although cigarette smoking accounts for the highest number of “attributable deaths” of any physical risk factor in the United States (13), relatively few mortality projections have taken this factor into consideration. Pampel (14) forecasts sex differences in mortality in high-income countries while using data on smoking. Based on the relationship between smoking prevalence and lagged smoking mortality, he predicts a rapid decline of the logged ratio of male-to-female smoking-related mortality between the 1990s and 2020 for an aggregate of 21 developed countries. Bongaarts (15) identifies mortality attributable to smoking and projects mortality assuming that component to be constant. Several national projection programs refer to smoking patterns as a basis for projecting narrowing sex differences in mortality (e.g., ref. 16).

The possibility that there is a strong cohort pattern in the relationship between smoking and mortality opens up promising new possibilities for mortality projection. A period-specific approach to smoking-based projections would require making predictions of smoking behavior at all ages for all future periods. In contrast, a cohort approach can rely primarily on cohort behaviors that have already been observed.

In this paper we incorporate a cohort smoking factor into the original Lee-Carter model to project mortality for the United States. We build on the Lee-Carter model by producing mortality forecasts jointly for males and females within the same model, rather than combining the two sexes or treating them individually. We use the average number of years spent as current smoker by age 40 as an indicator of smoking history for

each cohort. This variable was successfully used to predict mortality by Preston and Wang (2). This variable is an indicator of smoking duration rather than of smoking intensity, which cannot be recovered from the retrospective data available. Studies of lung cancer mortality find that duration is a much more successful predictor than intensity (17).

Methods and Data

Model. We assume that men’s and women’s mortality is influenced by age and period in the manner hypothesized by Lee and Carter, but that it is also influenced by smoking histories. In addition, we assume that males and females share a common temporal trend in mortality. The new model is:

$$\ln(M_{x,t}^{\text{Male}}) = a_x^{\text{Male}} + b_x^{\text{Male}} \cdot k_t + c \cdot S_{x,t}^{\text{M}} + \varepsilon_{x,t}^{\text{Male}} \quad [2.1]$$

$$\ln(M_{x,t}^{\text{Female}}) = a_x^{\text{Female}} + b_x^{\text{Female}} \cdot k_t + w \cdot S_{x,t}^{\text{F}} + \varepsilon_{x,t}^{\text{Female}} \quad [2.2]$$

In the above equations, a_x^{Female} and a_x^{Male} are the constant elements in the mortality profile by age for women and men, respectively; b_x^{Female} and b_x^{Male} describe the changes in the mortality rates at age x in response to changes in k_t over time for women and men, respectively; k_t is the level of mortality at time t , assumed to be the same for men and women when smoking histories are controlled; $S_{x,t}^{\text{F}}$ and $S_{x,t}^{\text{M}}$ are cohort smoking indices for women and men by single year and single age; and c and w are coefficients that indicate the effect of a particular smoking history on mortality of men and women, respectively. To reiterate, the differences between this modified Lee-Carter model and the original Lee-Carter model are that we introduce cohort smoking behaviors into the models and that men and women share a common mortality trend, k_t .

Eqs. 2.1 and 2.2 are then combined into a single equation for estimation purpose after a simple transformation as follows:

$$\ln(M_{x,t}^{\text{Male}}) - c \cdot S_{x,t}^{\text{M}} = a_x^{\text{Male}} + b_x^{\text{Male}} \cdot k_t + \varepsilon_{x,t}^{\text{Male}} \quad [2.1.1]$$

$$\ln(M_{x,t}^{\text{Female}}) - w \cdot S_{x,t}^{\text{F}} = a_x^{\text{Female}} + b_x^{\text{Female}} \cdot k_t + \varepsilon_{x,t}^{\text{Female}} \quad [2.2.1]$$

We can combine the above two equations into one:

$$\begin{bmatrix} \ln(M_{x,t}^{\text{Male}}) - c \cdot S_{x,t}^{\text{M}} \\ \ln(M_{x,t}^{\text{Female}}) - w \cdot S_{x,t}^{\text{F}} \end{bmatrix} = \begin{bmatrix} a_x^{\text{Male}} \\ a_x^{\text{Female}} \end{bmatrix} + \begin{bmatrix} b_x^{\text{Male}} \\ b_x^{\text{Female}} \end{bmatrix} \cdot k_t + \begin{bmatrix} \varepsilon_{x,t}^{\text{Male}} \\ \varepsilon_{x,t}^{\text{Female}} \end{bmatrix}. \quad [3]$$

We obtain Eq. 3 by simply stacking Eqs. 2.1.1 and 2.2.1 vertically. For example,

$$\begin{bmatrix} \ln(M_{x,t}^{\text{Male}}) - c \cdot S_{x,t}^{\text{M}} \\ \ln(M_{x,t}^{\text{Female}}) - w \cdot S_{x,t}^{\text{F}} \end{bmatrix}$$

is a new matrix in Eq. 3 obtained by stacking matrices $(\ln(M_{x,t}^{\text{Female}}) - w \cdot S_{x,t}^{\text{F}})$ and $(\ln(M_{x,t}^{\text{Male}}) - c \cdot S_{x,t}^{\text{M}})$ vertically.

Eq. 3 can be estimated by singular value decomposition. In particular, we find a set of parameters that minimizes the sum of the squared errors, E , which is defined in the following form:

$$E = \text{Sum} \left\{ \text{diag} \left(\left(\begin{bmatrix} \ln(M_{x,t}^{\text{Male}})' \\ \ln(M_{x,t}^{\text{Female}})' \end{bmatrix} - \begin{bmatrix} \ln(M_{x,t}^{\text{Male}}) \\ \ln(M_{x,t}^{\text{Female}}) \end{bmatrix} \right) \right. \right. \\ \left. \left. * \left(\begin{bmatrix} \ln(M_{x,t}^{\text{Male}})' \\ \ln(M_{x,t}^{\text{Female}})' \end{bmatrix} - \begin{bmatrix} \ln(M_{x,t}^{\text{Male}}) \\ \ln(M_{x,t}^{\text{Female}}) \end{bmatrix} \right)^T \right) \right\},$$

where

$$\begin{bmatrix} \ln(M_{x,t}^{\text{Male}})' \\ \ln(M_{x,t}^{\text{Female}})' \end{bmatrix}$$

is the estimated matrix given $c = c'$ and $w = w'$. E is equal to sum of all of the squared elements in the error matrix, which is

$$\left(\begin{bmatrix} \ln(M_{x,t}^{\text{Male}})' \\ \ln(M_{x,t}^{\text{Female}})' \end{bmatrix} - \begin{bmatrix} \ln(M_{x,t}^{\text{Male}}) \\ \ln(M_{x,t}^{\text{Female}}) \end{bmatrix} \right).$$

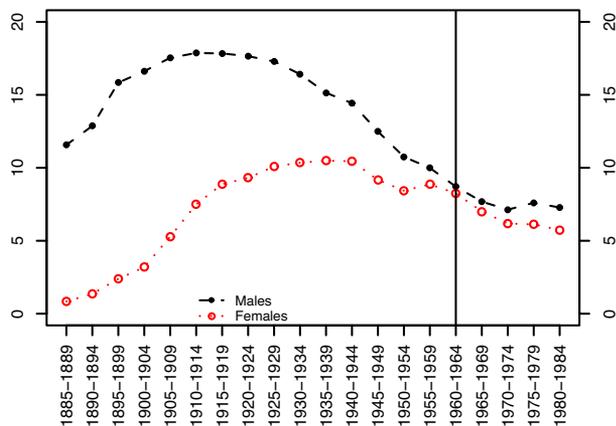


Fig. 2. Average number of years spent as a cigarette smoker before age 40 among men and women in different birth cohorts. Source: Preston S, Wang H (2).

Data. Mortality data are derived from the Human Mortality Database (www.mortality.org). The population data in this source are drawn from the U.S. Census Bureau. Data on deaths are taken from the National Center for Health Statistics. We use single-year of age mortality rates from ages 50 to 84 for every calendar year from 1969 to 2003. The first year for which we have complete smoking histories for all cohorts in the age interval under study is 1969. The cohort smoking index is reconstructed from the work by

Table 1. Estimated values of a_x and b_x for males and females

Age, yr	a_x^{Male}	a_x^{Female}	b_x^{Male}	b_x^{Female}
50	-5.4288	-5.7211	0.01233	0.01589
51	-5.3452	-5.6295	0.01157	0.01394
52	-5.2601	-5.5362	0.01263	0.01401
53	-5.1882	-5.4542	0.01287	0.01360
54	-5.1121	-5.3831	0.01457	0.01400
55	-5.0140	-5.2890	0.01437	0.01332
56	-4.9326	-5.1942	0.01451	0.01312
57	-4.8562	-5.1168	0.01544	0.01290
58	-4.7449	-5.0023	0.01459	0.01179
59	-4.6879	-4.9402	0.01546	0.01192
60	-4.5846	-4.8391	0.01595	0.01256
61	-4.5110	-4.7527	0.01461	0.01058
62	-4.4063	-4.6454	0.01514	0.01159
63	-4.3424	-4.5808	0.01522	0.01026
64	-4.2714	-4.5078	0.01518	0.01028
65	-4.1795	-4.4067	0.01596	0.01146
66	-4.1108	-4.3342	0.01527	0.01057
67	-4.0312	-4.2455	0.01571	0.01122
68	-3.9397	-4.1479	0.01538	0.01152
69	-3.8637	-4.0646	0.01568	0.01240
70	-3.7654	-3.9512	0.01666	0.01403
71	-3.6969	-3.8729	0.01503	0.01240
72	-3.5976	-3.7543	0.01637	0.01398
73	-3.5210	-3.6644	0.01577	0.01389
74	-3.4440	-3.5800	0.01623	0.01486
75	-3.3561	-3.4747	0.01680	0.01579
76	-3.2725	-3.3823	0.01633	0.01490
77	-3.1908	-3.2905	0.01599	0.01445
78	-3.1151	-3.1937	0.01544	0.01489
79	-3.0232	-3.0833	0.01586	0.01532
80	-2.9173	-2.9643	0.01680	0.01605
81	-2.8291	-2.8666	0.01550	0.01435
82	-2.7363	-2.7497	0.01614	0.01483
83	-2.6381	-2.6359	0.01582	0.01499
84	-2.5467	-2.5254	0.01618	0.01496

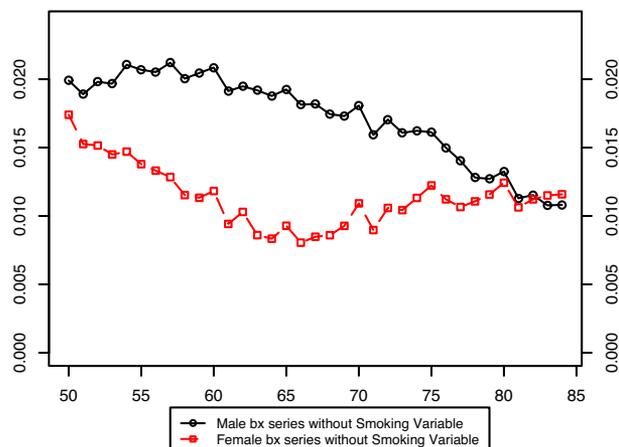


Fig. 3. Comparison of male and female b_x values in a model without smoking histories: United States, 1969-2003.

Burns *et al.* (ref. 7, with supplement supplied by the author), who relied on National Health Interview Surveys conducted between 1965 and 2001. We have converted the data into estimates of the average number of years that members of various birth cohorts smoked before age 40. We have assigned the smoking behavior for a 5-year birth cohort to each single-year birth cohort within that group. The value of the smoking variable for various birth cohorts of men and women is shown in Fig. 2.

Results

We solved simultaneously for all parameters, iterating across all possible values of c and w . The sum of squared errors was minimized by the set of parameters $c' = 0.031$ and $w' = 0.017$. A greater responsiveness of males to smoking is consistent with the epidemiologic studies reviewed above. In particular, Cancer Prevention Study I found that smoking raised male death rates by 91% and female rates by 46% (3), a sex difference very similar to that estimated here. Table 1 presents the estimated values of a_x^{Male} , a_x^{Female} , b_x^{Male} , and b_x^{Female} .

We have estimated Eq. 3 with and without data on smoking. The addition of the smoking variable explains 20% of the variance that was left unexplained by the Lee-Carter Model without smoking: the modified R^2 (i.e., the proportion of variance explained after age effects are accounted for) increases from 0.965 to 0.972. In addition, parameter values and projections are strongly affected.

Fig. 3 presents the estimated values of the b_x for males and

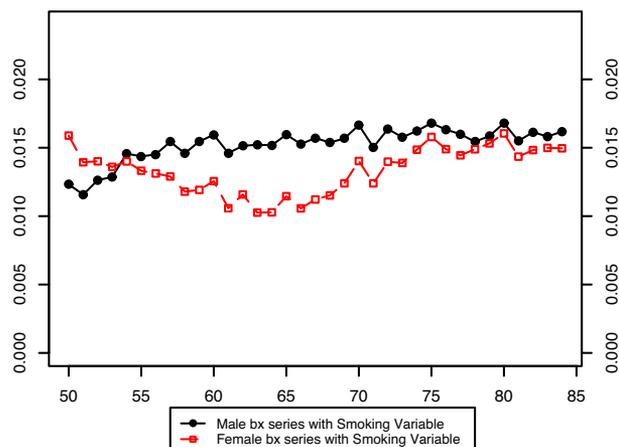


Fig. 4. Comparison of male and female b_x values in a model with smoking histories: United States, 1969-2003.

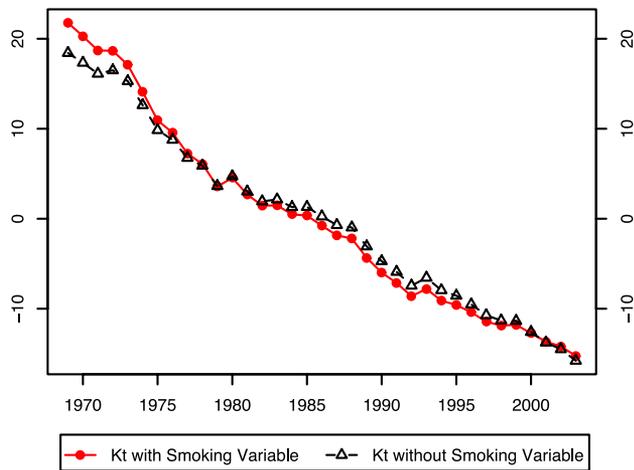


Fig. 5. Comparison of k_t values in models with and without smoking histories: United States, 1969–2003.

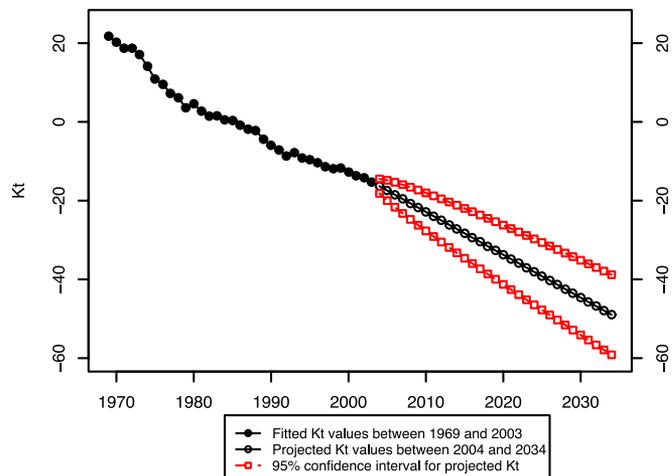


Fig. 6. Fitted and projected k_t values: United States, 1969–2034.

females without the smoking variable in the model. It is evident that there is a large disparity between the sexes in the estimated pattern of change in age-specific death rates when smoking is not included in the model. Furthermore, the age pattern of change is irregular for both sexes. Fig. 4 compares the b_x of males and females when the smoking variable is introduced into the model. Including smoking produces much more similarity in age patterns of mortality change between men and women. And when smoking is included in the model, both the male and female patterns of change in age-specific death rates are much more level with age than when smoking is omitted.[†]

This result suggests that changing smoking histories have distorted the observed age/sex pattern of change in death rates during this period. Younger ages experienced larger mortality declines than older ages in Fig. 3, probably because declines in cumulative smoking histories were sharper for younger ages than for older ages during the period of observation (Fig. 2). The fact that the male distortion in Fig. 3 relative to Fig. 4 is larger than the female distortion is also consistent with this interpretation, because male smoking has declined more sharply. Thus, accounting for smoking iron out many anomalies in the shape and sex differences in estimated b_x , anomalies that may have otherwise become exaggerated as the projection period advanced.

Fig. 5 shows the values of k_t produced by models with and without the smoking variable. Differences between the k_t series are not striking, but it should be recalled that in projections, k_t is multiplying the b_x , so that the same values of k_t may produce much different projections in the models with and without smoking variable.

Projections. The first step in the projection process is forecasting the mortality index, k_t . We followed the standard procedures suggested by Box and Jenkins (18) and Makridakis *et al.* (19) to find the appropriate ARIMA time series model for the mortality trend, k_t . ARIMA (0, 1, 0)—in other words, random walk with drift—best describes our index. Using this model, we forecasted the values of k_t to year 2034. We ended the projection period in 2034 because that is the latest year for which we have actual data

(in some cases, partial) on smoking behavior for all cohorts in the age range 50–84. Fig. 6 presents the fitted and projected values of k_t , with 95% confidence intervals for the projected values.

To forecast age-specific mortality rates between 2004 and 2034 for the age interval 50–84, we also required smoking indices for all cohorts who will be alive during this period. Most of those cohorts were already beyond age 40, and their smoking behavior could simply be observed. For the 4 cohorts born 1965–1969, 1970–1974, 1975–1979, and 1980–1984, we must project some portion of their smoking histories.

Our estimation strategy was to predict the smoking index (cumulative years smoked by age 40) based on cumulative years smoked at successively earlier ages by using the experience of cohorts who have already reached age 40 years. In other words, we predicted the value of smoking by age 40 years based on regressions with independent variables representing cumulative smoking indexes by age 35, by age 30, by age 25, and by age 20 years. We also added a trend variable to the regressions. Regressions are estimated on data for the 16 cohorts for which we have complete data up to age 40 years. The regressions in all cases explain at least 97% of the variance in cumulative years of smoking before age 40 years. We used these models to estimate

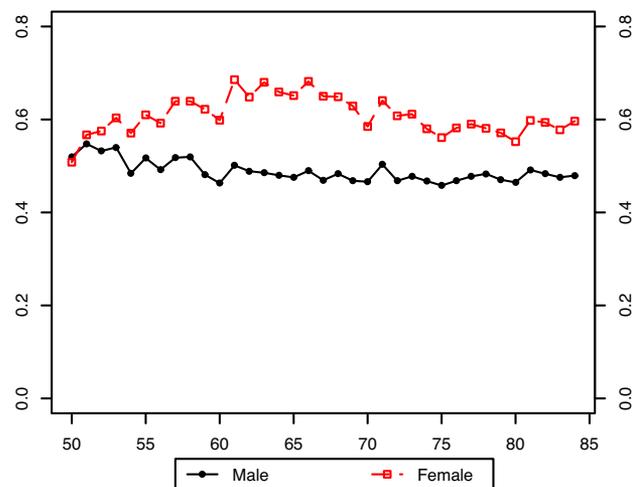


Fig. 7. Ratio of projected age-specific mortality rate in 2034 to age-specific mortality rate in 2003 using a model with smoking histories.

[†]Note that a level set of b_x would be produced by a family of Gompertz curves in which mortality change occurred through the “environmental” term. According to the well-known Gompertz curve, $\ln M_x = A + B \cdot X$, where A refers to environmental factors that are independent of age and B is the age slope of mortality. Two populations with the same value of B but different values of A would have age-specific death rates that were parallel on a log scale. That result would produce a level set of b_x in the Lee–Carter model. Thus, a level set of b_x is consistent with a common description of the age pattern of mortality, the Gompertz model, whereas an irregular or sloped series is not.

