

November 3, 2000
date last saved: 11/03/00 3:16 PM
date last printed: 11/21/00 4:00 PM

Evaluating the Performance of Lee-Carter Mortality Forecasts

Ronald Lee
Demography and Economics
University of California
2232 Piedmont Ave.
Berkeley, CA 94720
rlee@demog.berkeley.edu

Timothy Miller
Demography
University of California
2232 Piedmont Ave.
Berkeley, CA 94720
tmiller@demog.berkeley.edu

Research for this paper was funded by a grant from NIA, R37-AG11761. We thank John Wilmoth for making available mortality data for the US, France, Sweden, and Japan, through the Berkeley Mortality Data Base. We thank Statistics Canada and Francois Nault for making Canadian data available. John Wilmoth and Ken Wachter provided very useful suggestions for the analysis. This paper is available in electronic form at www.demog.berkeley.edu.

I. Introduction

Important policy decisions are made today based on forecasts of the elderly population far into the future. Public pension policies are the prime example. Fundamental changes have been proposed for the US Social Security system in part because of a projected financial crisis 37 years from now, driven by population aging. Old age dependency ratios are the key variable in these forecasts, and they depend on the number of elderly in the numerator, and the number of working age people in the denominator. The denominator depends heavily on future trends in fertility and perhaps migration, and these are notoriously difficult to forecast. The elderly in the numerator have already been born, at least for forecasts up to a 65 year horizon, and so it is on firmer ground. Yet Keilman (1997) finds systematic under-prediction of the elderly population in a number of industrial nations. He reports that after 15 years, forecast errors of –15% are not uncommon for elderly females 85+ (Keilman, 1997:272). A recent National Academy of Sciences study reports that UN projections done between 1965 and 1990 had average errors of about –10% for the elderly populations of Europe and North America after 15 years (net of baseline error; National Research Council, 2000:46; average errors in Third World countries were smaller but also negative). While immigration must have contributed to these errors, the main culprit is the systematic under-prediction of mortality decline and life expectancy gain. The National Academy study also reports that “For industrial countries, increases in life expectancy have been under-projected” (p.132) and Keilman (1998:38) reaches a similar conclusion. We will suggest that these problems continue in the recent and current forecasts of industrial nations, including those by the US Social Security Administration.

Mortality forecasts are typically based on the subjective judgments of the forecaster, sometimes buttressed by expert opinion, and it is these judgments that have tended to underestimate the pace of subsequent mortality decline in recent decades. Another approach, not without its own problems, is to reduce the role of judgment by using extrapolation to forecast mortality.

Recently, Lee and Carter (1992, henceforth LC) developed a new extrapolative method for modeling and forecasting mortality based on the analysis of long term trends, and used it to make probabilistic forecasts of US mortality to 2065. Since that time, the method has attracted a certain amount of attention and acceptance as well as a number of criticisms. The most recent Census Bureau population forecasts (Hollmann et al., 2000) use the Lee-Carter forecast as a benchmark for their long-run forecast of US life expectancy. The two most recent Social Security Technical Advisory Panels have recommended the adoption of the method, or forecasts consistent with it, by the Trustees. The method has also been used to forecast mortality in a number of other countries (most recently for the G7 nations, see Tuljapurkar et al., 2000).

Our main purpose in this paper is to make a careful and detailed assessment of the performance of the Lee-Carter method for forecasting mortality. The possibilities for the kind of ex post analysis of forecast performance by Keilman and the NAS panel, reported

above, are limited, but we will examine performance over the nine years which are available since the jump-off in 1989. However, because the method involves less subjective judgment than those that have been used in the past, it is possible to construct hypothetical forecasts with jump-off years earlier in the 20th century, pretending we had only the data available up to that point, and comparing the subsequent pseudo-forecasts to the actual outcomes. We also conduct some similar but less detailed experiments using the method to produce forecasts for Japan, Canada, France and Sweden, with jump-off year in 1950. Last, we examine age patterns of decline during the 20th century and consider the possibility that the age pattern has changed over time contrary to the assumptions of the method. We discuss the suggestions and criticisms that have been made in light of the results of these studies.

II. Overview of the LC approach

The basic LC model of age-specific death rates (ASDRs, and denoted $m_{x,t}$) is:

$$\ln(m_{x,t}) = a_x + b_x k_t + \varepsilon_{x,t} \quad (\text{Equation 1})$$

Here a_x describes the general age shape of the ASDRs, while k_t is an index of the general level of mortality. The b_x coefficients describe the tendency of mortality at age x to change when the general level of mortality (k_t) changes. When b_x is large for some x , then the death rate at age x varies a lot when the general level of mortality changes (as with $x=0$ for infant mortality, for example) and when b_x is small, then the death rate at that age varies little when the general level of mortality changes (as is often the case with mortality at older ages). Note that the model assumes that all the ASDRs move up or down together, although not necessarily by the same amounts, since all are driven by the same period index, k_t . In principle, not all the b_x need have the same sign, in which case movement in opposite directions could occur. In practice, all the b_x do have the same sign, at least when the model is fit over fairly long periods. Note that the proportional rate of decline of any death rate is give by $b_x (dk/dt)$. If dk/dt is constant, that is if k_t is declining linearly, then each ASDR will decline at its own age specific exponential rate, proportional to b_x , and depending on the rapidity of the decline in k_t . The same model was selected by Gomez de Leon (1990) using exploratory data analysis on the historical data for Norway, out of a larger set of possibilities.

The strategy is to estimate this model on the historical data for the population in question, obtaining values for a_x , b_x and k_t . The values of k_t form a time series, with one value for each year of data. Standard statistical methods can then be used to model and forecast this time series. LC selected a random walk with drift as the appropriate model, which has the form:

$$k_t = k_{t-1} + c + e_t \quad (\text{Equation 2})$$

In this specification, c is the drift term, and k is forecast to decline linearly with increments of c , while deviations from this path, e_t , are permanently incorporated in the trajectory. The variance of e_t is used to calculate the uncertainty in forecasting k over any given horizon. The drift term, c , is also estimated with uncertainty, and the standard error of its estimate can be used to form a more complete measure of the uncertainty in forecasting k .

The projected k can then be used in Equation 1, together with the estimated a_x and b_x , to calculate forecasts of the ASDRs, and from these any desired life table functions can be derived. The probability intervals on the forecasts of k can then be used in the same way to calculate intervals for the forecasts of the ASDRs, and (because these are all linear functions of the same k) the forecast of e_0 . However, forecast errors in the ASDRs and e_0 derive additionally from the $\varepsilon_{x,t}$ and from uncertainty about the true values of a_x and b_x . LC show that these latter sources of error matter less and less as the forecast horizon lengthens, and they are dominated by uncertainty about k in the long run. For a forecast horizon of 10 years, 98% of the standard error of the forecast of e_0 is accounted for by uncertainty in k ; for the individual age-specific rates, the other sources of uncertainty are more important initially and remain important longer, but after 25 years most account for less than 10% of the standard error of the forecasts (see LC Table B2).

From inspection of Equation 1 it is apparent that there is no observed variable on the right hand side of the equation, so ordinary regression methods cannot be used to estimate the model. LC describes a simple approximate method using regression methods, but the Singular Value Decomposition (SVD) gives an exact least squares fit. Also note that if a_x , b_x and k_t form one set of coefficients for the model, then a_x , b_x/A and $A*k_t$ will be an exactly equivalent set, for any constant A . Similarly, $a_x - b_x*A$, b_x , $k_t(1+A)$ will be an equivalent formulation for arbitrary constant A . LC stipulated a unique representation by setting a_x equal to the average of the logarithms of $m_{x,t}$ over the data period, and setting the average value of k_t equal to zero. In this case the sum of the b_x values is unity.

The method has a number of appealing features. The basic model is very simple, and although its use for forecasting involves a number of steps, each is simple in itself. The method is “relational” in demographers’ terminology. That is, it involves the transformation of actual existing mortality schedules for each study population, and therefore on the one hand is largely non-parametric, and on the other hand incorporates particular features of the mortality pattern of a given population. The method is also probabilistic, in the sense that it involves statistical fitting of models, and the quality of the fit of the historical data can be used to provide probability intervals for the forecasts. As a matter of empirical fact, in the applications of the method to date, involving at least ten national data sets, the historical trend in k has always been found to be highly linear with time, and the random walk with drift has been found to give a good fit. This approximate linearity is useful for forecasting. It contrasts with the typically nonlinear trajectories of life expectancy, which rises at a decelerating rate when age-specific

mortality rates decline at constant exponential rates. Nonetheless, it is clear that if a data series extends sufficiently far back in time, the linearity of decline would cease to hold. Finally, the method can also be used as the basis of a simple model life table system, and indirect estimation methods can be developed to expand the mortality data available as the basis for forecasting.

There have been a number of criticisms and suggestions since this original article was published, and the original method has been modified and extended in various ways. Some have thought that the probability bands are implausibly narrow (e.g. Alho, 1992:673). Others have argued that many age-specific rates are so low that they can't realistically be projected to decline much further. Some argue that biomedical information should inform the forecasts, perhaps through incorporating expert opinion as is done by the Social Security Actuaries. Some have called for more within-sample testing of the methods, and others have questioned whether the a_x and b_x should be treated as invariant. Bell (1997) noted that the model did not fit the jump-off data very well.

Considerable work has been done to refine and extend the method since the original LC article. Wilmoth (1993) has developed improved fitting methods based on weighted SVD. Methods for modeling and forecasting regional systems of mortality have been developed (Lee and Nault, 1993). Better procedures for dealing with the jump-off year have been developed (Bell, 1997). Alternatives for modeling mortality for the oldest old have been explored. Consideration has been given to the special role of leader and follower countries (Wilmoth, 1998). The method has been applied to cause of death data (Wilmoth, 1998) to sexes separately, and by race (Carter and Lee, 1992; Carter 1996). There have been many applications to countries other than the US (e.g., Lee and Rofman, 1994; Tuljapurkar et al., 2000). Lee (2000) provides a summary of the model's development, extensions, and applications such as stochastic forecasts of social security system finances.

III. Assessing the original forecast

In their original article, LC noted that the model would not fit the age-specific mortality data exactly in the jump-off year, which would mean that the initial conditions for the forecast would not be quite right. This would inevitably lead to error which would be particularly important in the early years of the forecast. They noted that it would be possible to set a_x equal to the most recently observed log age specific-rates, and thereby fit the initial conditions exactly (with $k_t = 0$). However, they argued that this practice might extrapolate idiosyncratic features of mortality in the jump-off year, and it was therefore preferable to estimate a_x as the average values of the log death rates (LC:665-666). In retrospect, this appears to have been a mistake, since the error in e_0 of .6 years at the jump-off year caused significant bias in the forecasts for the first decade, as we shall see below, and as Bell (1997) has pointed out (LC estimated e_0 for 1989 at 75.66 years, whereas official data puts it at 75.08). Bell (1997) assessed the performance of four mortality forecasts: LC (as published); LC (with the jump-off year corrected); McNown-

Rogers; and the Social Security actuaries (SSA). He concluded that the LC forecasts did better than the SSA or McNown-Rogers, but that a corrected LC forecast did better still.

Figure 1 displays the original LC mean forecast of e_0 , a similar forecast but with the correct jump-off level, and the SSA projections done at the same time. The bias in the original LC projections is apparent, but it is also apparent that those projections correctly identified the trend in e_0 . SSA appears to be somewhat low, ending up about 0.8 years below the actual e_0 . The adjusted LC projection is about 0.2 years too low in 1998 (the latest data available to us). Over this period, the actual e_0 always remains well within the 95% prediction interval for both the original LC and the adjusted LC.

If the forecasts of e_0 performed well from 1989 to 1998, how about the forecasts of the individual age-specific rates? Once again, there are certainly errors due to the errors in initial conditions. Figure 2 instead focuses on the LC projected age-specific rate of decline of death rates from 1989 to 1997 for sexes combined, since this will not be affected by the errors in initial rates. It also plots the actual rates of decline, and those projected by SSA. The agreement between the LC forecast and the actual rates of decline is striking, particularly at the older ages. The SSA projections, however, incorrectly forecast slower mortality decline in the young adult years. We will return to this topic later, for a different perspective on the age pattern of decline.

IV. Analysis of LC Projection Errors in Hypothetical Historical Projections

A. The nature of the tests

In the original LC article, there were some tests of forecast performance within the historical data period, but none of these involved re-estimating a_x , b_x and k_t . Tests were restricted to forecasting k_t from different starting points in the historical period.. Here we will make a more rigorous test, in which we completely refit the model on each chosen sub-sample of data. Our earliest experimental forecast is based on data from 1900 through 1920. Our next uses data 1900 through 1921; our next 1900 through 1922; and so on until our last forecast uses data from 1900 through 1997 to make a forecast for 1998. In this way, we have 78 different forecasts for mortality one year ahead; 77 for a two year horizon; and finally one with a 78 year horizon. We re-estimated the a_x and b_x for each set of data, and then re-estimated k_t for these years conditional on these a_x and b_x estimates, by choosing k_t (in the second stage) so as to match exactly the given value of e_0 in the data for that year.¹ This departs slightly from the procedure in the original LC, where k_t was chosen to match total deaths, which requires annual age-distributed population data as well.

Once k_t was estimated for each year of the sample, we did not carry out standard diagnostic methods to choose an optimal ARIMA model for each data sub-sample, but rather assumed that the random walk with drift model held. It was fitted and used to

forecast k_t over the desired time range, and probability distributions were derived from the subsample ARIMA errors.

LC introduced a dummy variable for the influenza epidemic of 1918. Our preference today is to include the dummy (permitting a one time positive change in k in 1918, followed by a one time equal negative change in k in 1919), and in the forecast to incorporate a $1/T$ chance of an identical positive and negative change in k occurring, where T is the length of the base period over which the model was fit. This has a small effect on both the mean and the variance of the forecast. We did not do this for these experimental forecasts, here described.

Although these retrospective tests are something like the ex post analysis of forecasting errors, there are also significant differences. First, the method was developed with the benefit of the preceding century of mortality experience, so it would be surprising if it failed to accord with it. Second, a forecaster would have to decide how far back in time to go in fitting the model to historical data. Mortality data for the US does not start until 1900, and then covers only for a limited number of the states. All our forecasts use data back to 1900, although our first forecast has a jump-off year of 1920. Third, we have assumed that a random walk with drift is the forecasting model always used, although the rate of drift is estimated afresh for each forecast. It would not be feasible to choose manually an optimal ARIMA model specification for each of the 78 forecast jump-off years. Had this been done, the short-run performance of the model would presumably have been better, but it is possible that the long-run performance would have been worse.

B. Forecasting to 1998 (e0)

Figure 3 plots all 78 forecasts for life expectancy in the year 1998, each from a different jump-off year, and each over a different forecast horizon. Each forecast for 1998 is plotted above its jump-off date. The 95% probability intervals are also plotted. The horizontal line indicates the observed value of life expectancy for 1998, so it is the true value relative to which the forecasts can be assessed. There are several points to note. *First*, although the hypothetical forecasts tend to be too low, they are generally fairly close to the actual value for 1998. The earlier forecasts, using data up through the 1920s and 1930s, are on average five years below the true value; beginning in 1946 all forecasts are within two years of the correct value. Over all, the mean forecasts look quite good. *Second*, the 95% probability intervals failed to contain the true value for 1998 in 12 out of the 78 forecasts, or 15% of the time, compared to the 5% which was intended. *Third*, the median forecast for 1998 fell below the actual value for 1998 in 74 of the 78 forecasts, or 95% of the time. That suggests downward bias.

C. Errors by forecast horizon (e0)

It is also useful to assess forecast errors (defined as forecast value – actual value) by horizon. We have done this for horizons of 1, 5, 10, 20, 40 and 60 years. For a one year horizon, we have 78 different jump-off dates, while for the 60 year horizon, we have only 19. For each forecast, we find the percentile in its probability distribution where the observed value falls. For example, if the actual value corresponds to the median of the forecast distribution, we assign it 50. If it corresponds to the lower 7% of the distribution,

we assign it 7; and so on. We then plot the frequency distribution of these percentile scores. If the probability distribution associated with each forecast does in fact describe the probability distribution of errors, then this frequency distribution should be uniform between 0 and 100. If the actual distribution of percentiles is more concentrated in the middle, around 50, that indicates that the distribution of the errors is more tightly clustered than our forecast leads us to expect, and if there are less in the middle of the distribution and more towards the 0 and 100 end, then our forecast understates the width of the error distribution. If most of the true values fall below the 50th percentile, then most of the time we have over-predicted, while if they fall above the 50th percentile, then we have tended to under-predict systematically the true value.

Figure 4 plots the histogram of the percentiles for each horizon. We can see that the actual forecast errors match the predictive distribution quite well for forecast horizons up to 10 years. By 20 years, a tilt towards positive errors is unmistakable, and this tilt intensifies at the 40 year horizon and again at 60 years. Note that the vertical scales are increasingly compressed as the errors become more concentrated.

Table 1 presents various measures of forecast performance, including the Mean Squared Error (MSE), the Mean Absolute Percent Error (MAPE), the average error (Bias), the percent of positive errors, and the proportion of actual values that fall within the 95% probability interval of the forecast. The table reports performance by forecast horizons as well as a summary over all forecast horizons.

Table 1

Forecast Horizon (N)	Average error	MAD	RMSE	MAPE	% under- projected	% within 95% interval
1-5 (380)	-0.11	0.45	0.60	0.16	54	99
6-10 (355)	-0.32	0.82	1.03	0.47	56	100
11-20 (635)	-0.73	1.23	1.60	1.15	67	97
21-30 (535)	-1.37	1.47	1.99	2.03	84	100
31-40 (435)	-1.68	1.73	2.14	2.45	91	100
41-50 (335)	-2.23	2.25	2.75	3.41	96	95
51-60 (235)	-3.54	3.54	3.75	5.07	100	89
61-78 (171)	-4.38	4.38	4.53	5.39	100	80
ALL (3,081)	-1.49	1.76	2.34	2.45	78%	97%

The average errors are negative, indicating that the method tended to under-predict gains in life expectancy in the US, particularly when launched from earlier dates. The percent under-predicted column confirms this. For example, 91% of errors for 31-40 year projection horizons were negative (predicted e_0 less than actual) and 100% of errors beyond a 50 year horizon were negative. The 95% confidence bounds contain the actual $e(0)$ value 97% of the time for all horizons combined. However, they appear to be too broad for intervals up to a 40 year horizon and too narrow for those beyond a 50 year horizon.

D. Error correlations by age, horizon

As noted briefly above, Equation 1 has an error term, $\varepsilon_{x,t}$, since the expression does not provide a perfect representation of variation in age-specific rates over time. In formulating the probability intervals for the life expectancy forecasts, this error term was ignored, and only errors arising from the innovation in k_t and from errors in estimating the drift term, were incorporated. If we were interested only in e_0 and if the $\varepsilon_{x,t}$ were uncorrelated across age, this assumption might be relatively harmless, because some twenty different values of $\varepsilon_{x,t}$ enter into the calculation of any life expectancy, and these will tend to cancel, leaving a small net effect. However, if the errors are correlated, such that those for older ages tend to move together and those for younger ages tend to move together, then they might have an important influence even on life expectancy. There are also errors in the estimation of the a_x and b_x coefficients, which are not taken into account in our probability intervals for the e_0 forecasts.

We find that forecast errors tend to be strongly positively correlated at younger ages, less so at older ages, and the errors at young ages are only weakly correlated with those at older ages. At longer horizons, correlations become more positive due to dominance of errors in k .

V. Assessing LC on historical time series from other countries

We also carried out a simple within-sample test for Sweden, Japan, France and Canada. For Canada, France and Sweden we constructed a forecast to 1995 from a jump-off date of 1950. For France and Sweden, we used data from 1900 to 1950. For Canada, data are available only from 1922 to 1950. For Japan, suitable data are available from 1950, so we took the later jump-off year of 1973. For France, we used dummies to capture the profound effects of both WWI and WWII but did not allow for a possible recurrence in the future, which would have greatly increased the variance of the forecast. Such decisions reflect the judgment of the analyst.

The results are shown in the panels of Figure 5. If the method had been used to forecast 1995 e_0 for Sweden, starting in 1950, it would have been right on target until 1980, and two years too low in 1995. Results for France and Canada are very similar. For Japan, forecasts from 1973 to 1996 are below the actual value, and are one year too low by 1996. Looking at all the forecasts combined, the 95% probability bounds contain the actual $e(0)$ values for 152 out of 162 forecasted values or 94% of the time, which is very

close to the 95% coverage predicted by the models. However, inspection shows that in every country there is a systematic tendency to under-predict future gains in life expectancy, just as in the US. We will return to this topic later.

VI. Changing age-shape of mortality

A number of people have suggested that the b_x coefficients might vary over time; this possibility was not explored by LC. Kannisto et al. (1994) found that the rate of mortality decline had been accelerating over recent decades for ages 80 to 100. Horiuchi and Wilmoth (1995) show that in a number of countries, mortality declines at older ages now take place more rapidly than at lower ages, reversing the historical pattern. This research suggests that it is important to take very seriously the possibility that the age pattern of mortality decline may change over time, and may not be well described by a fixed set of b_x coefficients. Note that the a_x coefficients will always be changing over different historical periods, because they are the average log death rates, and these averages will change in level as mortality falls, and change in shape because the b_x coefficients tell us that at different ages, mortality declines at different rates. This poses no problem, because the changing shape and level of the a_x are implicit in the b_x , and no additional treatment is necessary.

Recall that our earlier examination of the post-publication performance of LC showed that it correctly forecast the age pattern of mortality decline as well as the increase in e_0 over the past 9 years. This suggests that the fixed b_x assumption has worked well.

However, a closer examination of the age pattern of decline in the US shows otherwise. Figure 6 plots the average rate of decline for sexes combined mortality by age for 1900 to 1950 and for 1950 to 1995. It suggests that there has been an important change, with mortality now declining at roughly the same rate across all ages above 15, whereas for the first half of the century it declined far more rapidly at the younger ages. Figure 7, which shows changes in the historical age pattern of mortality decline in Japan, Sweden, Canada, and France, indicates similarly striking alterations, with a flattening of the age profile of decline.

Is this a long term change, rooted in the changing cause structure of mortality, or in the resistance of mortality at different ages to biomedical progress? Or is it due to what we might hope will be more transitory influences on young adult mortality in industrial nations, such as AIDS and accidents? We are not sure. But the more prudent course is to assume that these changes are long term, and to incorporate them into our forecasts in one way or another. A simple and satisfactory solution, adopted by Tuljapurkar et al. (2000), is to base the forecast on data since 1950, and assume fixed b_x over that range but not over the whole century. Only about 6% of life table deaths now occur in the age range affected by the changing age pattern, say from 10 to 50, so the changing age pattern of decline has relatively weak effects on the forecast of e_0 . It seems likely that the systematic tendency of the LC method to under-predict gains in e_0 at long horizons is in some way due to this changing age pattern of decline, but it is not clear exactly how.

VII. Comparison of official forecasts from SSA and others to LC forecasts

A. Forecasting to 1998

We have examined the historical record of SSA projections, including two earlier ones that were used by SSA but prepared by other agencies. Figures 8 and 9 examine forecasts of e_0 for the year 1998. Figure 9 compares the middle series forecast from SSA with the median LC forecast. The figure shows that the official projections have been systematically too low – by 12 years in 1930, by about 7 years in the 1940s, and then by 2 to 4 years until those done in 1980. In 1980, the SSA forecasts jumped too high for a few years, then dropped down too low again. It can be seen that the SSA estimates reacted strongly to the slow mortality gains of the 1960s, and then to the rapid gains of the 1970s. By contrast, the LC method responds only modestly to these fluctuations, since they only modestly affect the average trend over the century. The LC method also tends to be somewhat low in early years, but performs substantially better than SSA. It would have been closer to the true value in 1998 for most forecast horizons. It picks up the correct track for 1998 considerably earlier.

Figure 9 shows the high-low range of SSA projections along with the 95% probability interval of LC. The true value of e_0 for 1998 lies beyond the high bound for most of the SSA forecasts up until 1970.

B. Errors by horizon, comparison to LC

In assessing errors by forecast horizon, we have restricted our sample to post-1950 government forecasts. We have only three early government forecasts (pre-1950) – which provided $e(0)$ forecasts for only a few select years in the future. This makes the analysis of errors by length of horizon complicated for these groups. For comparison to LC, we use both the full sample (1920-1997) and a restricted sample which matches the time period of the SSA forecasts (1950-1997). For LC, we have hypothetical forecasts for every year. For SSA, the forecasts have been issued irregularly until 1980, after which they are available annually. In our calculations we have weighted each SSA forecast by the reciprocal of the number of forecasts issued within the decade. In this way, each decade contributes equally to the error estimates. Without weighting, the SSA results are dominated by forecasts done since 1980, and the longer horizon forecasts count for little.

Figure 10 compares the average bias in the SSA and LC forecasts by length of forecast horizon. Horizons are by single year from 1 to 7 and then grouped (8-12, 13-17, 18-22, 23-27, 28-38, 39-46, and 39-60 years). SSA forecasts issued since 1950 compare favorably with LC forecasts issued since 1920 for horizons up to 15 years, and do worse thereafter. However, when we compare only those LC forecasts issued during the same time period (since 1950) as the SSA projections, we find that LC performs very much better at all horizons.

Figure 11 compares the root mean square error (RMSE) for SSA and LC forecasts, again weighting each SSA forecast by the reciprocal of the number of forecasts issued within the decade. Once again the SSA forecasts since 1957 do better than the LC forecasts from

1920. However, when we compare the LC forecasts made over the post-1950 time period, then at all horizons beyond two years, the LC perform better than SSA and substantially better as the forecast horizon increases.

C. General problem of official forecasts

Long-run government forecasts have relied on expert opinion which proved to be too pessimistic about the future. This pessimistic outlook might be attributed to the mood of the country at the time the forecasts were issued. Two of the earliest population forecasts were produced by the National Resource Committee (1937) during the Great Depression and by the National Resource Planning Board (1943) during the second World War. And yet, at those times, the data were telling a different story, since mortality had been declining quite rapidly over the previous decades. A quote from the 1943 report is telling in this regard. Thompson and Whelpton state their objection to statistical forecasting methods such as extrapolation: “More important, the extrapolation of past trends according to such formulas might show future trends which seemed incompatible with present knowledge regarding the causes of death and the means of controlling them.” (National Resources Planning Board, 1943, p. 10). This suggests an alternative explanation for the pessimistic bias of expert opinion: present knowledge informs us about current limits, but not the future means of overcoming them. The Lee-Carter approach bases its long-run forecasts on the century-long decline in mortality in which limits have been continuously confronted and overcome.

VIII. Conclusions

We can extract the following lessons from these investigations:

- 1) The LC forecasts of life expectancy and the age pattern of mortality performed quite well for the period since publication, at least after adjusting for an error in jump-off level.
- 2) Hypothetical LC projections from various historical jump-off dates in the 20th century would have preformed well. For forecasts with jump-off after 1945, LC projections are always within two years of the actual e_0 in 1998. However, the forecasts tend to under-predict future gains, especially those in the distant future. Although the 95% probability bounds contain the true value of e_0 97% of the time, the bounds appear to be too broad for horizons up to 40 years and too narrow for those beyond 50 years.
- 3) Social security projections also have systematically under-predicted gains in e_0 since 1950. The average error and mean squared error for LC forecasts since 1950 are substantially lower than those of SSA since 1950, when each decade is given equal weight.
- 4) LC life expectancy forecasts for Canada, Sweden and France with jump off year 1950 and for Japan with jump-off year 1973 would have performed very well. However, as in the US, the forecasts would have systematically under-predicted actual gains, particularly at longer horizons.

- 5) Contrary to a basic assumption in the Lee-Carter model, the age pattern of mortality decline has shifted systematically in the US, Sweden, France, Canada, and Japan in the second half of the 20th century, with a flattening of the age-specific rates of decline above age 15. While this has distorted e_0 projections only slightly, it can have a substantial effect on projected death rates at ages from 1 to 35. For example, the median forecast for the US life expectancy in 2075 based on post-1950 data is 86.2 years- about 0.5 years higher than the forecast based on post-1900 data. The age-specific rates for ages 1 to 35 are 30% to 80% lower in the forecast based on post-1900 data. However, the absolute errors are small because the projected rates themselves are so low.
- 6) Overall, the results suggest that the LC method produces surprisingly good forecasts over rather long time periods. Used for long-term forecasts within the 20th century, it would generally have tended to under-predict future gains. The probability intervals, despite some problems, also do a surprisingly good job of containing the true outcomes.

These findings bear on some of the criticisms and suggestions addressed to the LC method, that were briefly mentioned earlier. Some suggested that the probability bounds were too narrow. We found that for forecasts up to 50 years into the future, the probability bounds are too broad rather than too narrow, but for longer forecasts they are somewhat too narrow with 80 to 90% coverage rather than the intended 95%. Some have argued that many age specific-rates are so low that they can't realistically be projected to decline much further. Although death rates at ages 10 to 50 continue to decline, their rates of decline have decelerated relative to those at other ages. These changes in the age distribution of decline may reflect an approach to lower limits. However, the declines at the younger and older ages continue unabated or have accelerated. Some have suggested that biomedical information should inform the forecasts. But if so, how? The Social Security Actuaries have used expert opinion on mortality decline by cause of death. We find that their forecasts have been systematically too low, more so than those of the LC extrapolative approach, and the mean squared errors of their forecasts have been greater than those of LC as well. Some have questioned whether the relative pace of decline by age should be treated as invariant over the century. They are correct. In the second half of the century mortality at older ages has declined more rapidly relative to that at younger ages than in the first half of the century. Some have suggested that it would be better to take the most recently observed age-specific death rates as the jump-off point for the forecasts, rather than the fitted age distribution in the jump-off year. Based on the experience of the past ten years, this point also appears to be correct.

The Lee-Carter method takes a simple extrapolative approach. It is easy to think of reasons why its long-run forecasts should fail. Indeed, we have uncovered a number of shortcomings in the method's performance. What impresses us overall, however, is not these shortcomings we have so far found, but rather that they are not larger and more numerous. On these tests, the method performs better than we had reason to expect, both in predicting the future (and pseudo future), and in indicating uncertainty.

These results suggest that projections using the Lee-Carter method should be taken seriously. For example, Tuljapurkar et al (2000:792) use this method to project for a number of industrial nations that e_0 will be 1 to 4 years higher in 2050 than indicated by official projections, with larger discrepancies for Japan. It may well be that the systematic under-prediction of life expectancy by national and international agencies is continuing today. As industrial nations strive to confront the long-term funding problems of their public pension systems, it is particularly important that they have realistic projections of mortality. While we cannot know what future mortality trends will be, we suggest that Lee-Carter type forecasts provide a useful baseline forecast for planning.

References

- Alho, Juha M. (1992) "Modeling and Forecasting the Time Series of U.S. Mortality," *Journal of American Statistical Association*, v.87 n.419 (September) p.673.
- Bell, William R. (1997) "Comparing and Assessing Time Series Methods for Forecasting Age Specific Demographic Rates" *Journal of Official Statistics* 13:279-303.
- Carter, Lawrence (1996) "Long-Run Relationships in Differential U.S. Mortality Forecasts by Race and Gender: Non-Cointegrated Time Series Comparisons", 1996 Annual Meetings of the Population Association of America, New Orleans, May 9-11, 1996.
- Carter, Lawrence and Ronald D. Lee (1992) "Modeling and Forecasting U.S. Mortality: Differentials in Life Expectancy by Sex," in Dennis Ahlburg and Kenneth Land, eds, *Population Forecasting*, a Special Issue of the *International Journal of Forecasting*, v.8, n.3 (November) pp.393-412.
- Gomez de Leon, Jose (1990) Empirical DEA Models to Fit and Project Time Series of Age-Specific Mortality Rates," Unpublished manuscript of the Central Bureau of Statistics, Norway (July).
- Hollmann, Frederick W.; Mulder, Tammany J.; and Kallan, Jeffrey E. (2000) "Methodology and assumptions for the population projections of the United States: 1999 to 2100," Population Division Working Paper No. 38, U.S. Bureau of the Census.
- Horiuchi, Shiro and John R. Wilmoth (1995) "The aging of mortality decline." Presented at the Annual Meeting of the Population Association of America, San Francisco, April 6-8, 1995.
- Kannisto, Vaino; Lauritsen, Jens; Thatcher, A. Rodger; and Vaupel, James W. (1994) "Reductions in mortality at advanced ages: Several decades of evidence from 27 countries." *Population and Development Review*, Vol. 20, No 4, pp. 793-810.
- Keilman, Nico (1997) "Ex-post errors in official population forecasts in industrialized countries." *Journal of Official Statistics* (Statistics Sweden) 13(3): 245-277.
- Keilman, Nico (1998) "How accurate are the United Nations world population projections?" *Population and Development Review* 24 (supplement): 15-41.
- Lee, Ronald D. (2000) "The Lee-Carter Method for Forecasting Mortality, with Various Extensions and Applications," *North American Actuarial Journal*, Vol. 4, No. 1, pp. 80-91.

- Lee, Ronald D. and Lawrence Carter (1992) "Modeling and Forecasting the Time Series of U.S. Mortality," *Journal of the American Statistical Association* v.87 n.419 (September) pp.659-671.
- Lee, Ronald D. and Francois Nault (1993) "Modeling and Forecasting Provincial Mortality in Canada," paper presented at the World Congress of the International Union for the Scientific Study of Population, Montreal, 1993.
- Lee, Ronald D. and Rafael Rofman (1994) "Modelacion y Proyeccion de la Mortalidad en Chile," NOTAS 22, no 59, pp. 182-213. Also available in English from the authors, titled "Modeling and Forecasting Mortality in Chile."
- National Research Council (2000) *Beyond Six Billion: Forecasting the World's Population*. Panel on Population Projections. John Bongaarts and Rodolfo A. Butatao, eds., Committee on Population, Commission on Behavioral and Social Sciences and Education. Washington, D.C.: National Academy Press.
- National Resources Committee (1937) Population Statistics. Material prepared for a study of population problems. Washington D.C.: U.S. Government Printing Offices.
- National Resource Planning Board (1943) Estimates of the future population of the United States, 1940-2000. Washinton D.C.: U.S. Government Printing Office.
- Tuljapurkar, Shripad; Nan Li and Carl Boe (2000) "A universal pattern of mortality decline in the G-7 countries" *Nature* 405 (June): 789-792.
- Wilmoth, John R. (1993) "Computational Methods for Fitting and Extrapolating the Lee-Carter Model of Mortality Change," Technical Report, Department of Demography, University of California, Berkeley.
- Wilmoth, John R. (1998) "Is the pace of Japanese mortality decline converging toward international trends?" *Population and Development Review* 24(3): 593-600.

ⁱ In all cases, the data we use are taken from the SSA data base, as maintained on the Berkeley Mortality Data Base web site, www.demog.berkeley.edu. The original LC article used NCHS data, and for the period before 1933 estimated age specific mortality and e_0 indirectly using the age distribution of the total population, total deaths per year, and the a_x and b_x coefficients as estimated by SVD from the data 1933 to 1987, after the death registration area was complete.

Figure 1: Forecasts of life expectancy from 1989.

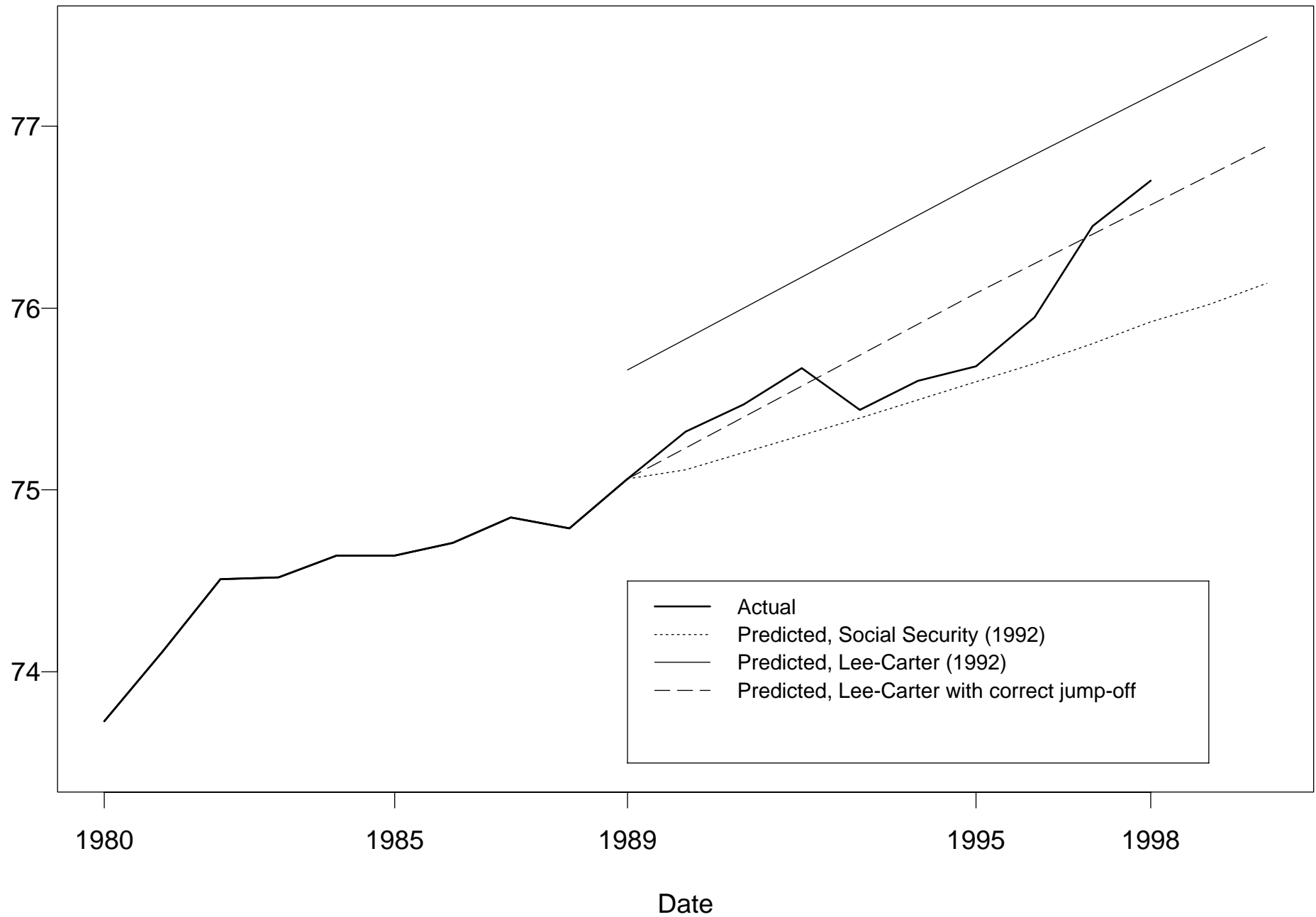


Figure 2: Average Annual Decline in Age-Specific Mortality, 1989-1997
Actual and Forecasts of Lee-Carter (1992) and SSA (1992)

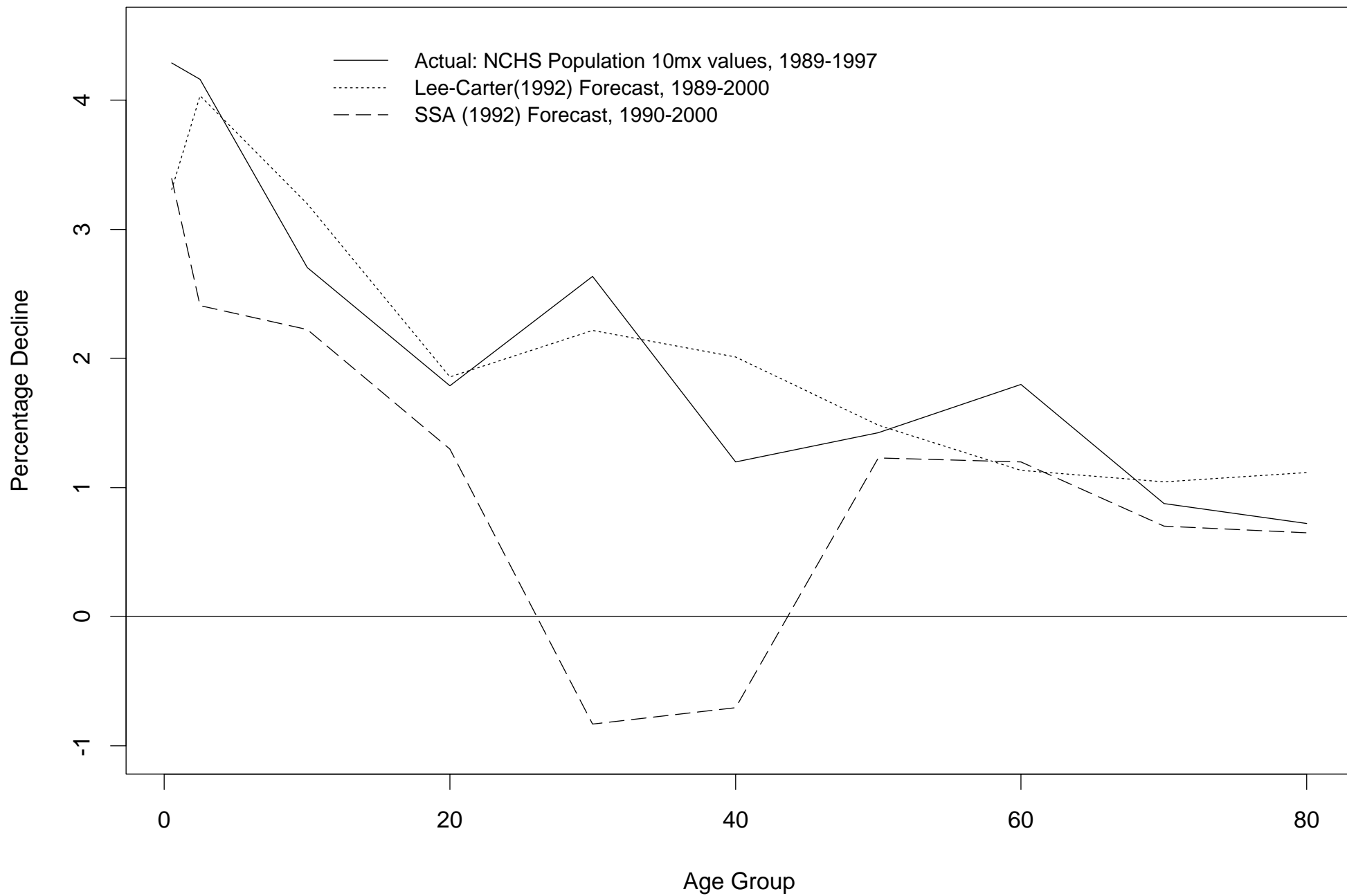


Figure 3: $e(0)$ Forecasts for the Year 1998 by Forecast Date

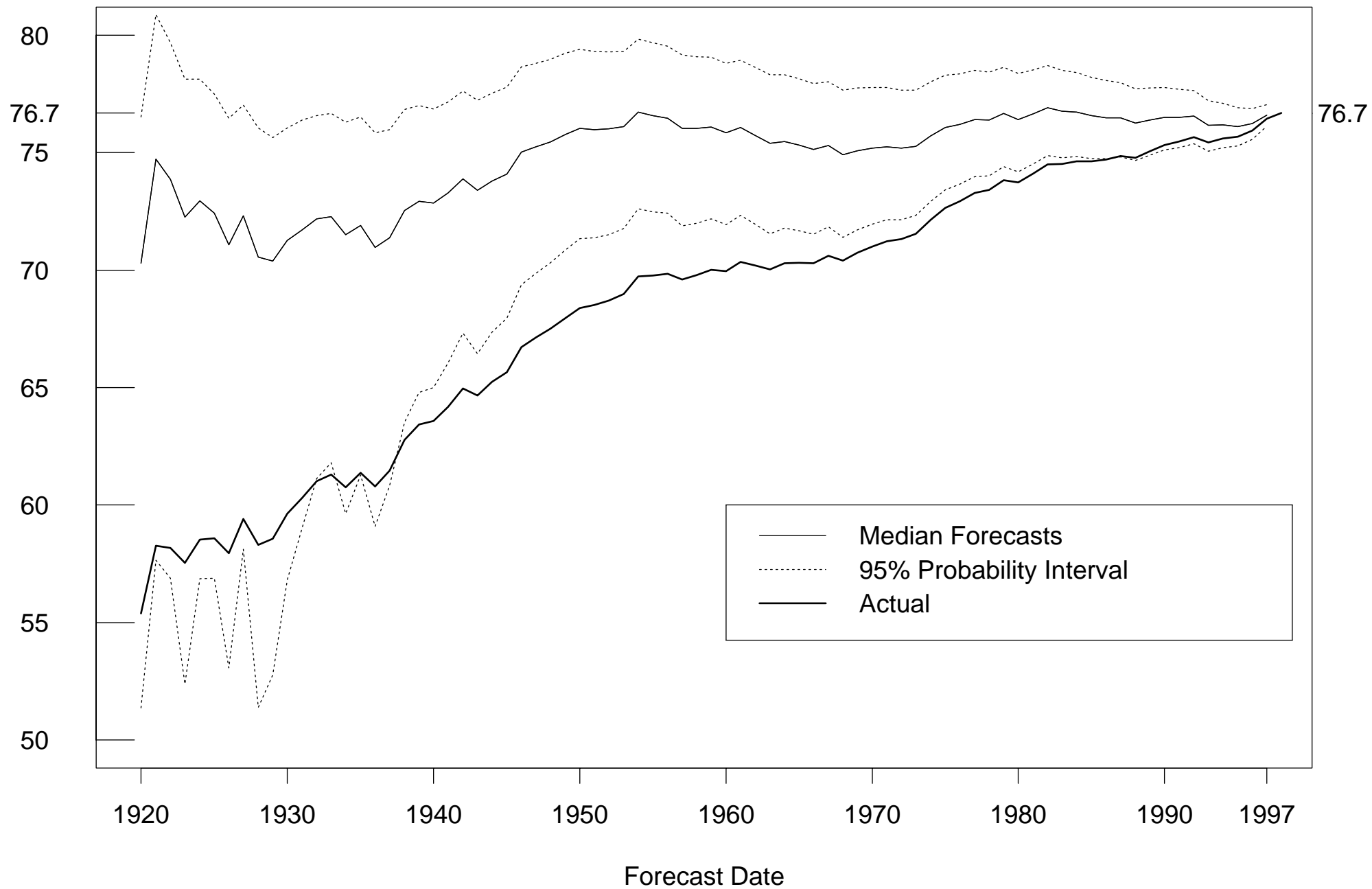
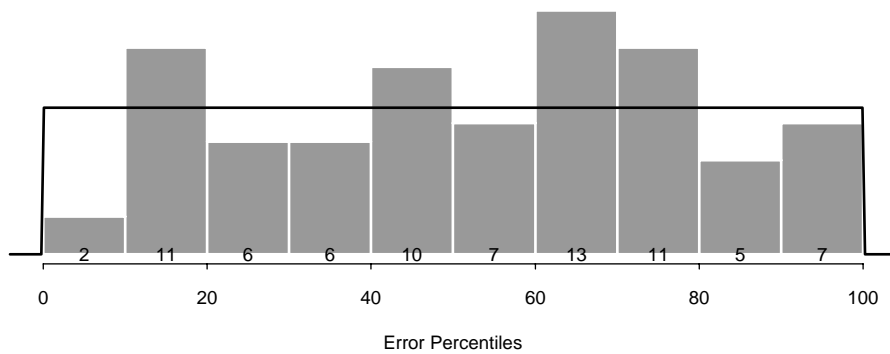
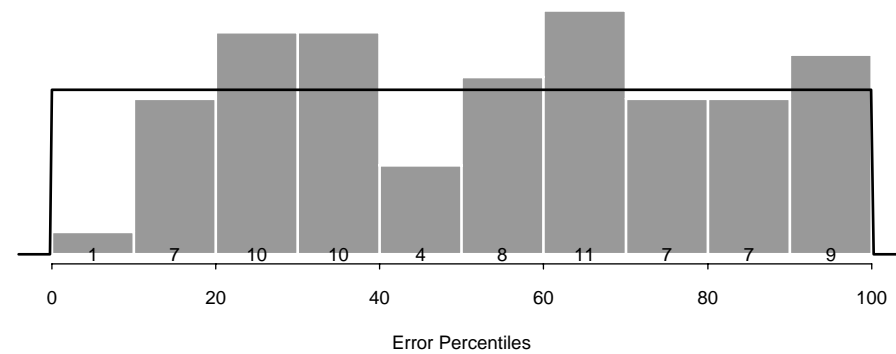


Figure 4: Percentile Error Distribution by Forecast Length

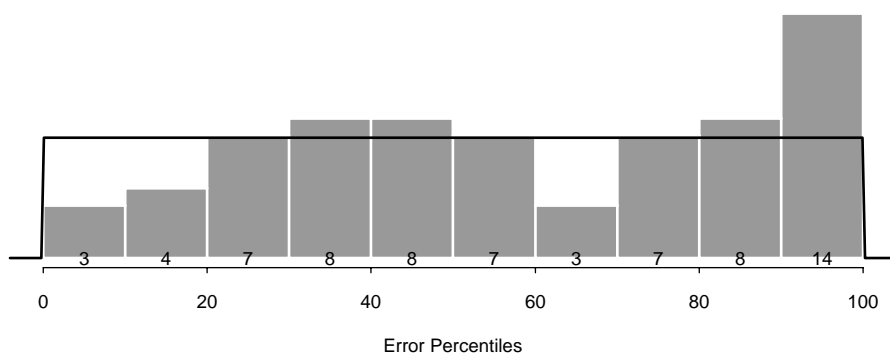
1 year forecasts, (78 obs.)



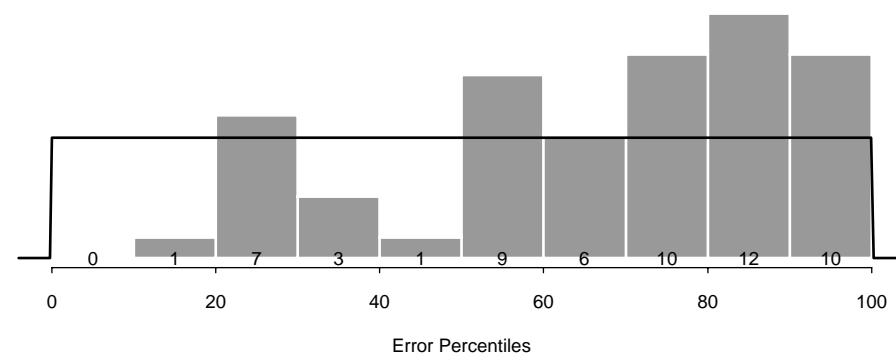
5 year forecasts, (74 obs.)



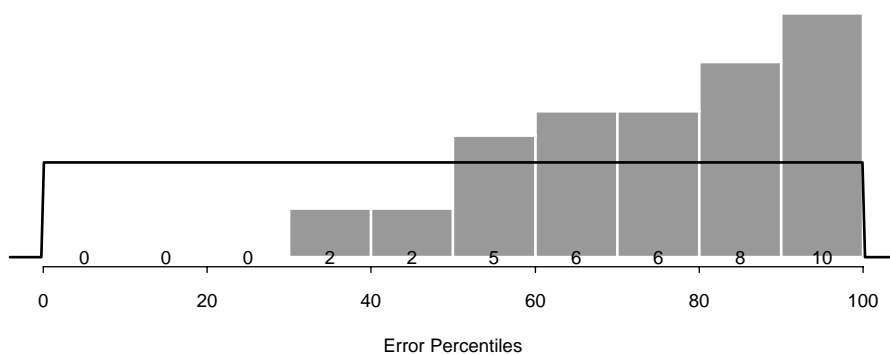
10 year forecasts, (69 obs.)



20 year forecasts, (59 obs.)



40 year forecasts, (39 obs.)



60 year forecasts, (19 obs.)

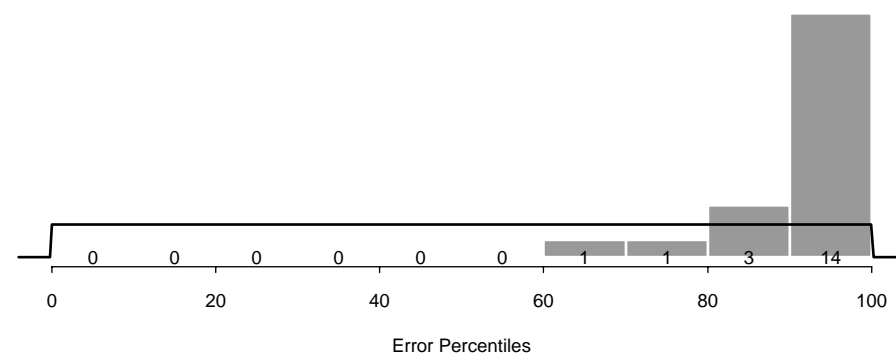
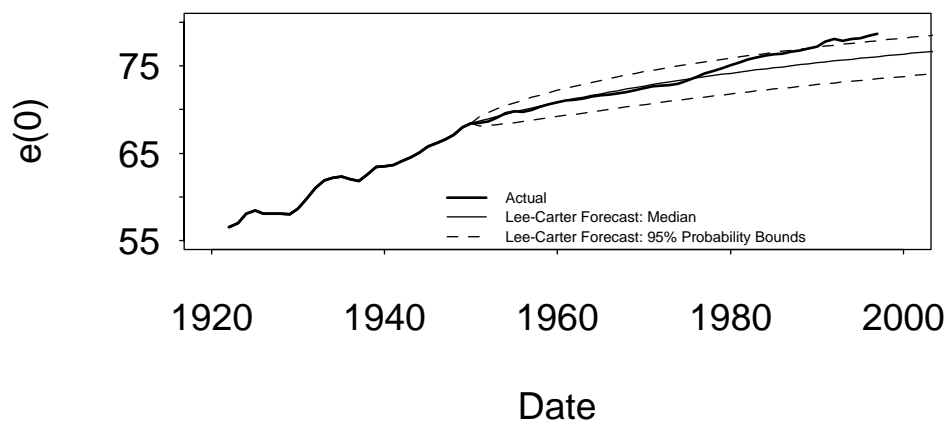
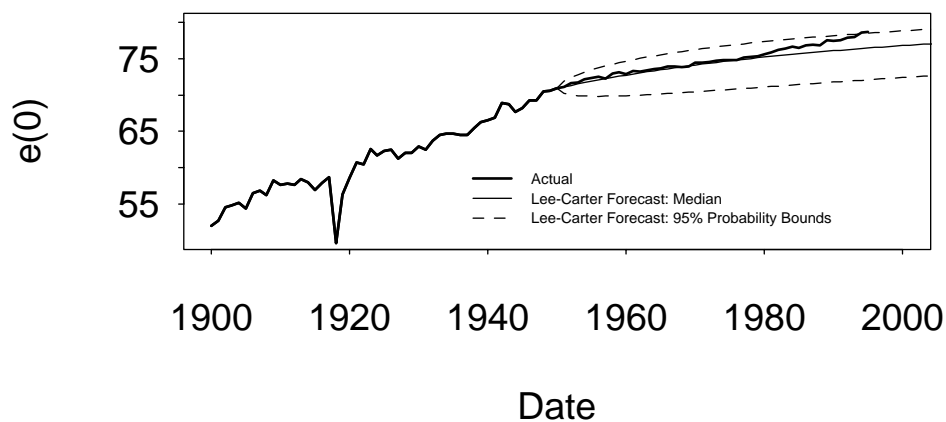


Figure 5: LC forecasts of life expectancy

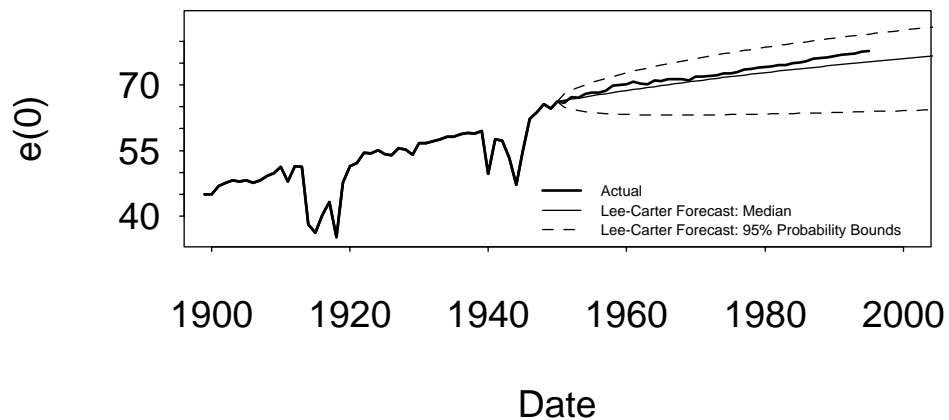
Canada from 1950



Sweden from 1950



France from 1950



Japan from 1973

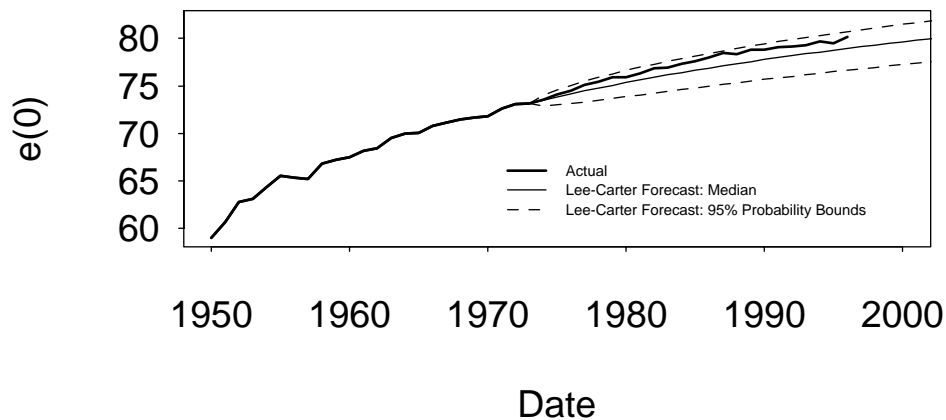


Figure 6: Average Annual Reduction in Age-Specific Death Rates, US

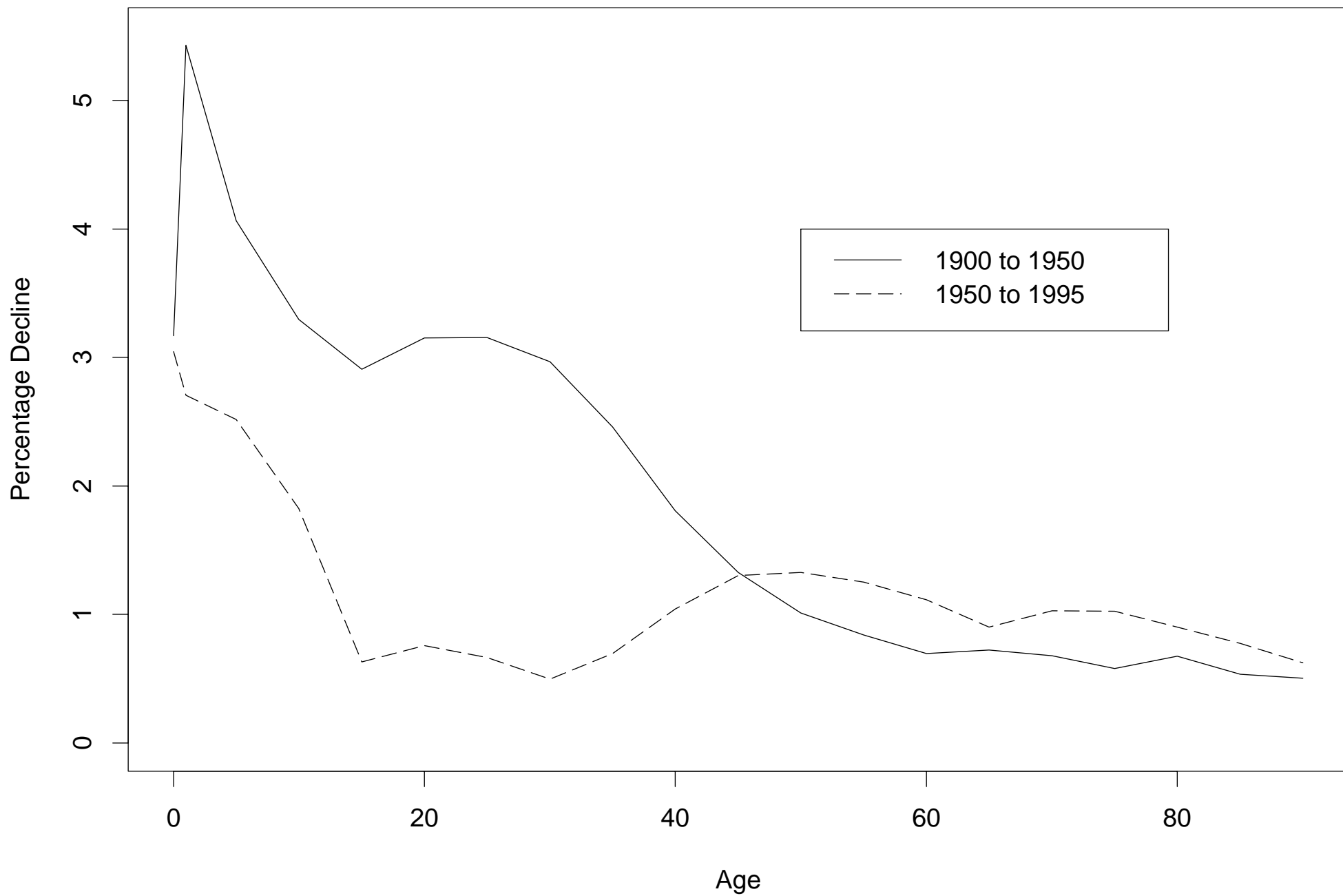


Figure 7: Average Annual Reduction in Age-Specific Death Rates

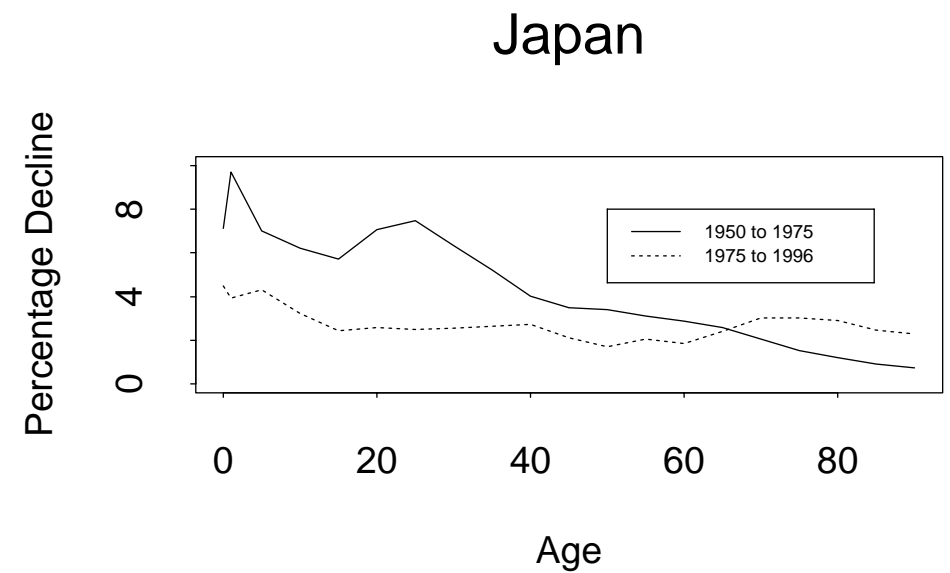
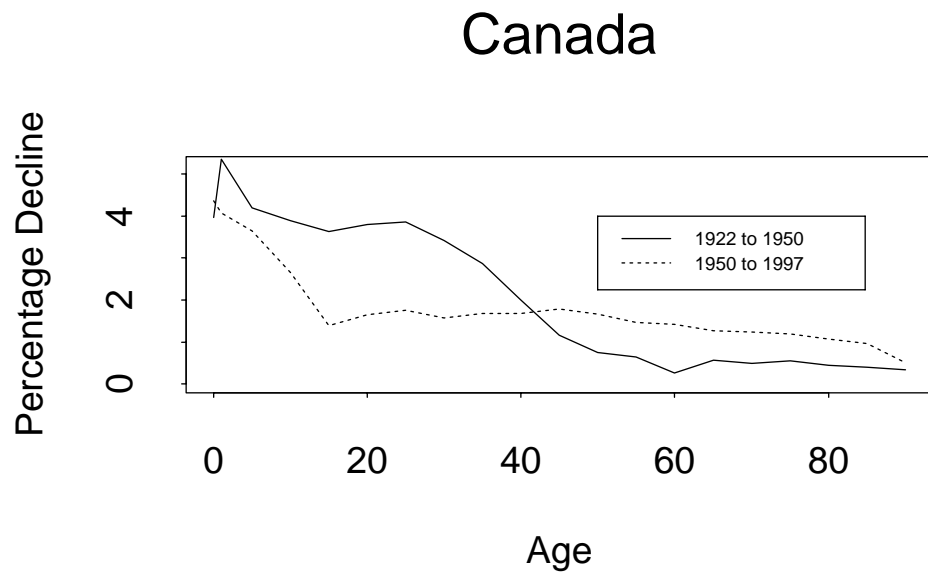
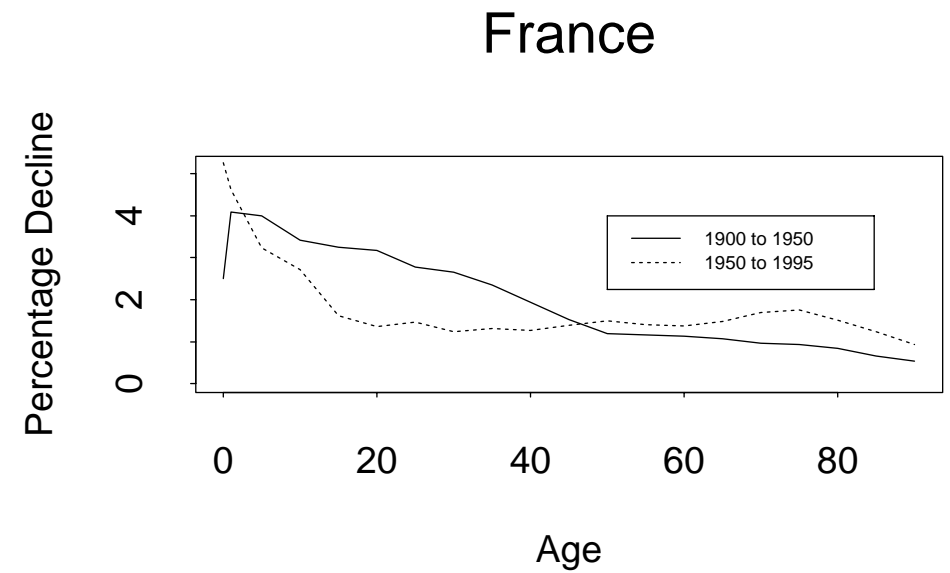
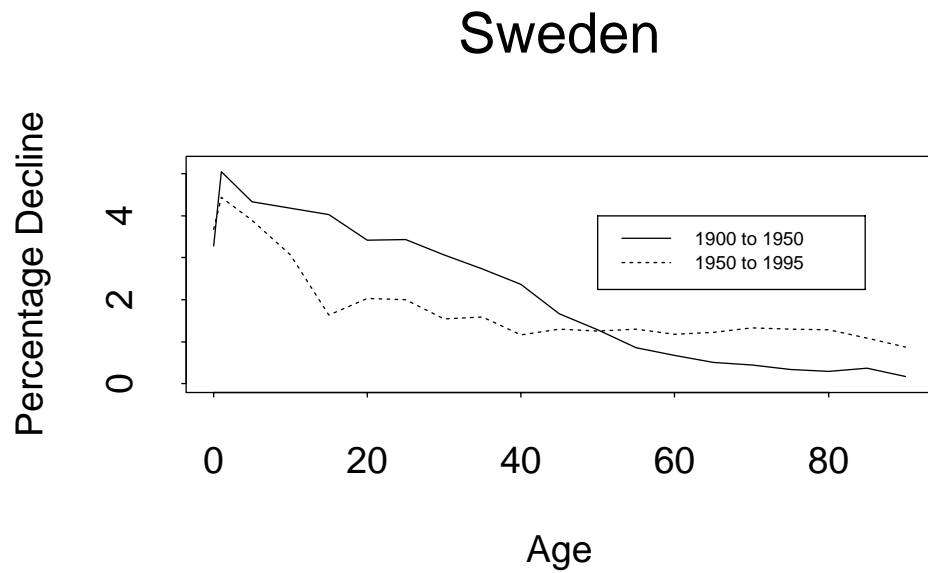


Figure 8: LC and SSA $e(0)$ Forecast for 1998, by Forecast Date

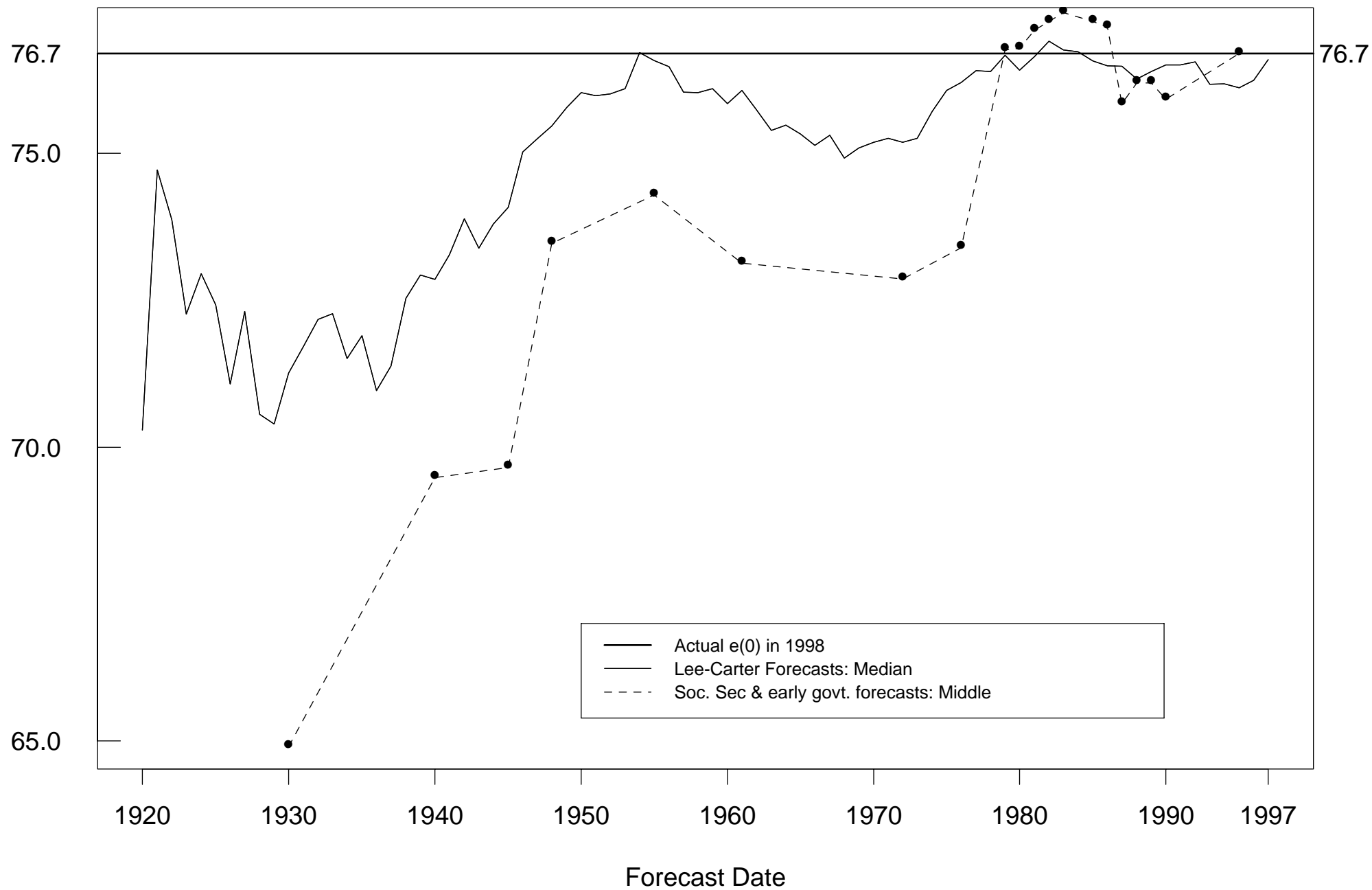


Figure 9: 95% Probability Interval and High-Low Range by Forecast Date

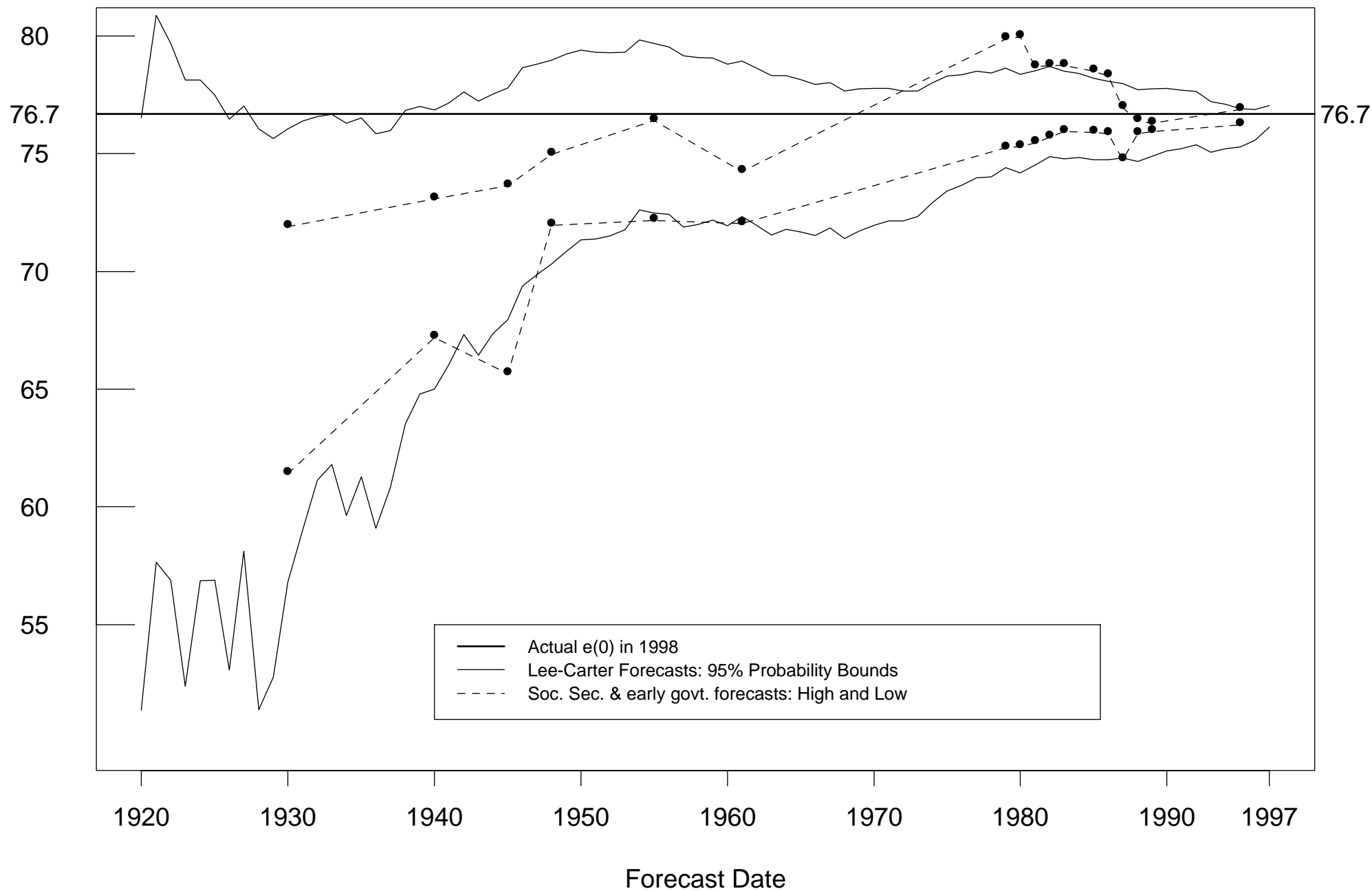


Figure 10: Average Bias in Forecasts of Life Expectancy

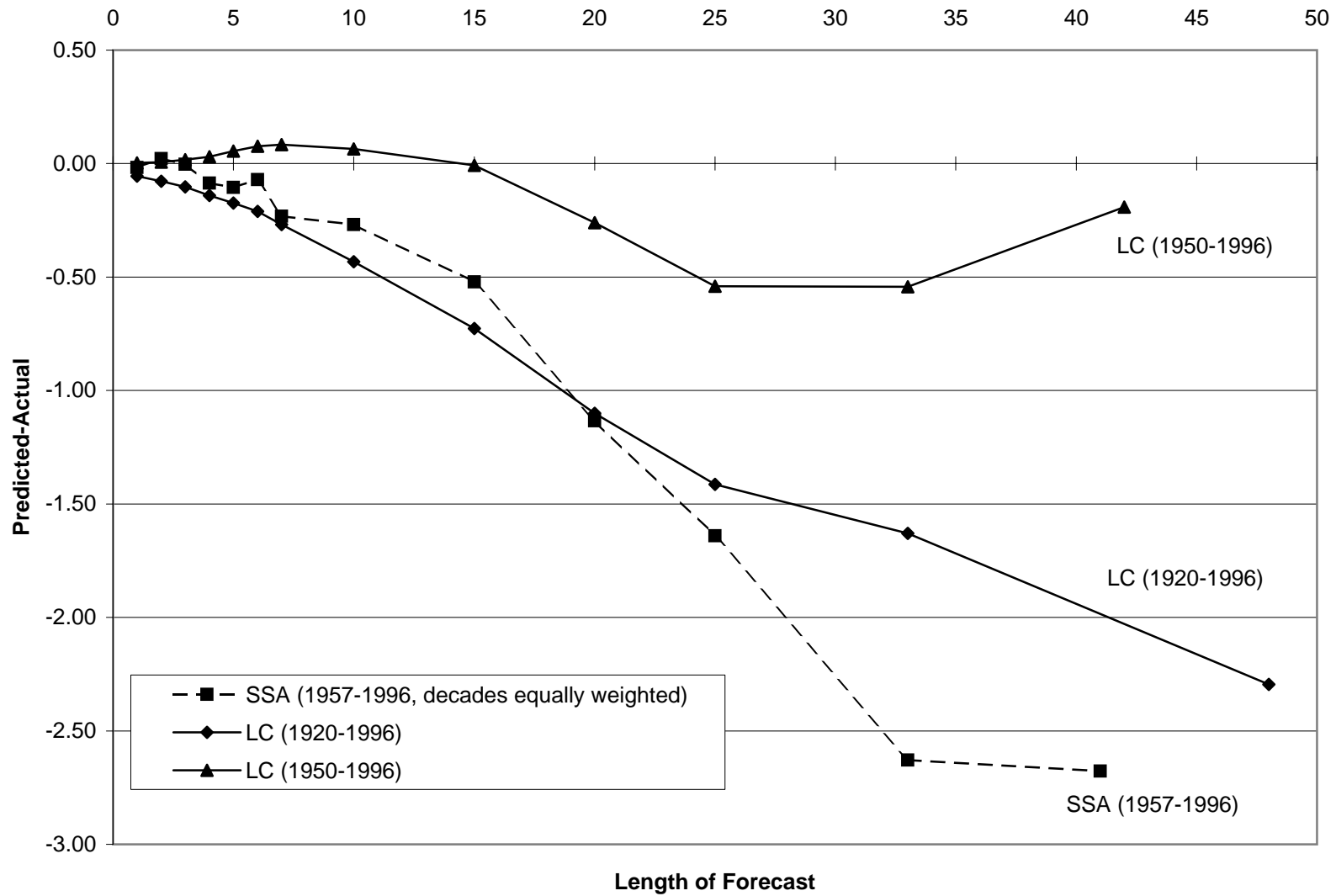


Figure 11: Root Mean Squared Error in Forecasts of Life Expectancy

