

Comparison of different approaches to incidence prediction based on simple interpolation techniques

Tadeusz Dyba^{1,2,*} and Timo Hakulinen^{1,3}

¹ Finnish Cancer Registry, FIN-00170, Helsinki, Finland

² Unit of Cancer Epidemiology, Karolinska Institute, S-17176 Stockholm, Sweden

³ Department of Public Health, FIN-00014, University of Helsinki, P.O. Box 41 (Mannerheimintie 172), Finland

SUMMARY

The paper compares three different methods for performing disease incidence prediction based on simple interpolation techniques. The first method assumes that the age-period specific numbers of observed cases follow a Poisson distribution and the other two methods assume a normal distribution for the incidence rates. The main emphasis of the paper is on assessing the reliability of the three methods. For this purpose, *ex post* predictions produced by each method are checked for different cancer sites using data from the Cancer Control Region of Turku in Finland. In addition, the behaviour of the estimators of predicted expected values and prediction intervals, crucial for investigation of the reliability of prediction, are assessed using a simulation study. The prediction method making use of the Poisson assumption appeared to be the most reliable of the three approaches. The simulation study found that the estimator of the length of the prediction interval produced by this method has the smallest coverage error and is the most precise. Copyright © 2000 John Wiley & Sons, Ltd.

1. INTRODUCTION

The prediction of disease incidence, including cancer incidence, is important as it forms the basis for the planning and organization of prevention, diagnosis and treatment in the community [1–3]. These activities are usually expensive. Thus, the reliability of the predictions used as the basis for undertaking practical measures is crucial to the decision process. Also, in a scientific context, verification of whether the prediction actually came true or not, which is the starting point for finding out the reasons for a possible discrepancy, calls for a measure of reliability [4]. From a statistical point of view, this measure, called the prediction interval, is the likely range of a future observation under the assumption that the model that has been chosen is correct.

*Correspondence to: Tadeusz Dyba, Finnish Cancer Registry, FIN-00170 Helsinki, Finland

†E-mail: tadek.dyba@cancer.fi

Contract/grant sponsor: Swedish Cancer Foundation

Contract/grant sponsor: Stockholm Cancer Society

Contract/grant sponsor: Ministry of Education in Finland

Contract/grant sponsor: Academy of Finland

The reliability of prediction directly depends on the choice of the model and the variables included in the model. Cancer incidence prediction is hindered by the lack of reliable knowledge of cancer aetiology, and also well recognized aetiological factors, for example, smoking, cannot be directly accounted for due to insufficient data. Thus, those factors must often be represented by surrogate variables such as age, period of observation, and year of birth (cohort).

A variety of alternative, often complex, models, even without theoretical support concerning their form, can fit the data within the range of observations (base of prediction). Alas, most of them fail to save this property for predicting outside this range. A good fit does not automatically guarantee a good prediction. Moreover, complex models, such as polynomial age–period–cohort models [5] with complicated mathematical form or many parameters, can produce very long prediction intervals [6] of no practical use. In this case, the most plausible and successful way of describing the future course of cancer incidence is to use simple linear models on an arithmetic or logarithmic scale [7]. The principle of parsimony is crucial in the context of prediction. However, one should remember that linear models assume no changes of underlying trends in the base of prediction and if such a change occurs further investigation will be necessary by applying alternative predictive models [8].

This paper presents a comparison of three predictive models based on simple extrapolation techniques, focusing on the reliability of the different methods. The main emphasis is on comparisons between a method based directly on the numbers of observed cases as Poisson variables and methods modelling age-adjusted and age-specific incidence rates as normally distributed variables. The comparison is drawn by checking *ex post* predictions produced by different approaches using real data from the Cancer Control Region of Turku, Finland. In addition, the behaviour of estimators of predicted expected values and prediction intervals for different models is analysed based on a simulation study.

2. MATERIALS AND METHODS

2.1. Calculating prediction intervals

Age is the most common risk factor in chronic diseases and information about it is almost always available. One of the simplest prediction techniques is to adjust cancer incidence for age and then perform the prediction based on a regression model for the rates [9]. A period specific age-adjusted incidence rate represents a single summary index for all age groups together. This measure has been widely used in analysing and comparing cancer incidence trends [10].

Directly standardized incidence rates are weighted sums of the age-period specific numbers of cases:

$$M_t = \sum_i (w_i/n_{it})c_{it}$$

where c_{it} and n_{it} are, respectively, the age-period specific numbers of cases and person-years for period t and w_i is the weight of age group i in the reference population used for the adjustment, with $\sum_i w_i = 1$. It is commonly assumed in epidemiology that the numbers of cases observed in each age-period specific group are independent and have Poisson distributions [11]. On the other hand, the statistical analysis is performed assuming that the weighted sums are approximately normally distributed [12]. For the purpose of prediction using this approach it is important that the projected population for the period of prediction is chosen to be the standard population used

for adjustment [9]. Two predictive models making use of the age-adjusted rate have been used [9]. For cancers with increasing incidence a model of the form:

$$E(M_t) = \alpha + \beta t \quad (1)$$

where t is a numerical variable for period, has been used with the parameters α and β . Its linear form controls against explosively increasing predictions. The second model is

$$\log(E(M_t)) = \alpha + \beta t \quad (2)$$

with a logarithmic scale employed for decreasing trends in order to avoid negative predicted rates.

This method cannot produce age-specific predictions and the standardization can potentially mask differences (for example, increases and decreases) between age-specific trends [13], which can bias the overall result. Also, the assumption of normality may be inadequate if the numbers of events in the age groups are small.

To avoid the first of these disadvantages, each age group could be modelled separately. This can be done by applying the following models:

$$E(c_{it}/n_{it}) = \alpha_i + \beta_i t \quad (3)$$

$$\log(E(c_{it}/n_{it})) = \alpha_i + \beta_i t \quad (4)$$

and

$$\log(E(c_{it}/n_{it})) = \alpha_i + \beta t \quad (5)$$

Practical experience has shown that models of the form

$$E(c_{it}/n_{it}) = \alpha_i + \beta t$$

have little value, as the assumption of common slope parameters across all age groups is too strong in diseases such as cancer, where incidence increases rapidly with age.

Using an assumption of normally distributed incidence rates, c_{it}/n_{it} , a prediction for increasing trends can be made based on model (3), whereas models (4) and (5) should be applied for diseases with decreasing incidence rates. As the variances of random components will differ between age strata, incidence rates in different age groups must be modelled separately. Model (5) leads to multivariate analysis because different age groups cannot be independently modelled due to sharing the same parameter β .

Ways of dealing with the prediction interval problem for normally distributed age-adjusted and age-specific rates are to be found in textbooks of statistics [14, 15]. Appendix A shows details of how to calculate predictions for model (4).

A method [16] has been proposed for calculating approximate confidence and prediction intervals for linear and log-linear models, both for the total numbers of cases and for the age-adjusted incidence rates. According to common practice in epidemiology, this method assumes that the age- and period-specific numbers of incident cases, c_{it} , are independent and follow Poisson distributions. Models (3), (4) and (5) have been considered. Although the functional form of these models is unaffected by the change of the error structure from normal to Poisson, the models become essentially different and have different log-likelihood functions. For clarity, the formulae (3), (4) and (5) with Poisson errors are subsequently labelled (3P), (4P) and (5P), respectively.

A model of the form:

$$E(c_{it}/n_{it}) = \alpha_i(1 + \beta t) \quad (6P)$$

can also be applied. Provided that it fits, it gives a smoother set of age-specific predictions and narrower age-specific prediction intervals than model (3P) [17]. However, the results of prediction for the total number of cases, which is the main concern of this paper, based on models (3P) and (6P), if both models fit, are approximately the same. Models (3)–(5) and (3P)–(5P) share in the context of prediction the useful non-identifiability property [17, 18] which guarantees that the age–period models can at the same time be regarded as age–cohort models. Special macros [19] are required to implement models (3P)–(6P) as the final results cannot be obtained routinely from any statistical package.

2.2. Empirical study

To check how the methods presented work in practice, *ex post* predictions were made using each approach. The number of new female cancer cases in the Cancer Control Region of Turku, Finland [20], in 1989–1993 (female population 377 000 in 1994) was predicted based on incidence data from 1954–1978. The data set, including numbers of observed cases and population or person-year data counts, were tabulated in five-year periods (1954–1958, 1959–1963, ..., 1974–1978), and five-year age groups (30–34, 35–39, ..., 85+) for the following common sites of cancer: breast; colon; corpus uteri; ovary; stomach; pancreas; rectum; lung; skin melanoma, and bladder. The final number of age groups used for a site-specific prediction depended on the site. The population in the period 1989–1993 was used in all incidence calculations and also as a standard for the direct adjustment for age.

2.3. Simulation study

The crucial part of the analysis of reliability of a prediction with concomitant prediction interval is the investigation of the coverage error of the prediction interval, that is, the absolute difference between the coverage probability and the nominal level of the prediction interval. The behaviour of estimators produced by a model, those of predicted expected value and a prediction interval is also essential. Therefore, a simulation study was performed.

For each age-period specific group, a simulated observed number of cases from a Poisson distribution was obtained. This was done for all periods that formed the estimation basis of the prediction and for the period for which the prediction was made. It was assumed that the pattern of incidence generating the data was linear over time and independent between age groups. The parameters of the Poisson distributions for the estimation basis of the prediction were obtained from fitting the linear model (3P) to the breast cancer data, which is a cancer with increasing incidence. The same model was also used to produce age-specific predicted expectations for the period of prediction. There were 12 age groups and five periods in the estimation basis of the prediction plus one period for which the prediction was done. The periods of the estimation basis and that for the prediction correspond to those in the empirical study, 1954–1978 and 1989–1993, respectively. Thus 72 observations were produced for each simulation. Models (1), (3) and (3P) were fitted to the simulated data set and the results from each fit were stored. The procedure was performed 5000 times, enabling assessment of the coverage probabilities of prediction based on each model and examination of the distributions of the estimators.

Table I. Observed numbers of female cancer cases in the Cancer Control Region of Turku, Finland, in 1989–1993 and the predicted numbers with 95 per cent prediction intervals produced by different models.

Site	Model	Predicted number	95 per cent prediction interval	Observed number
Breast	(1)	1862	1559–2165	2095
	(3)	1862	1694–2029	
	(3P)	1840	1685–1996	
Colon	(1)	513	370–655	532
	(3)	513	428–597	
	(3P)	508	415–602	
Corpus uteri	(1)	458	380–535	428
	(3)	458	397–518	
	(3P)	455	375–535	
Ovary	(1)	475	325–625	413
	(3)	475	396–554	
	(3P)	473	400–559	
Stomach	(2)	404	353–462	386
	(4)	427	179–676	
	(4P)	429	351–508	
Pancreas	(1)	407	201–612	342
	(3)	407	292–521	
	(3P)	412	335–490	
Rectum	(1)	365	207–523	310
	(3)	365	293–437	
	(3P)	364	288–440	
Lung	(1)	214	40–386	285
	(3)	214	129–298	
	(3P)	214	139–289	
Skin melanoma	(1)	161	126–196	209
	(3)	161	119–203	
	(3P)	164	120–208	
Bladder	(1)	122	79–164	147
	(3)	122	86–157	
	(3P)	126	80–173	

The same procedure was applied for models (2), (4) and (4P), in order to check the coverage probabilities and the characteristics of the estimators of models making use of logarithmic transformation. Model (4P) fitted to data on stomach cancer was chosen as the model generating the data.

3. RESULTS

3.1. Empirical study

All models fitted well to the empirical data with one exception (Table I). The fitting of model (3P) for lung cancer indicated that the data were overdispersed so the prediction interval was adjusted

Table II. Empirical characteristics of the estimator of the width of the prediction interval for models used in the analyses, based on simulations of 5000 observations.

Model	Coverage probability	Coverage error	Mean value	Standard deviation	Minimum value	Maximum value
(a)						
Model (1)	0.8432	0.1068	256.7	107.7	17.1	803.1
Model (3)	0.9235	0.0265	276.3	38.5	138.7	409.2
Model (3P)	0.9498	0.0002	289.7	3.1	275.3	301.0
(b)						
Model (2)	0.7932	0.1568	135.0	62.6	3.0	424.1
Model (4)	0.9521	0.0021	109.2	338.4	38.2	14382.5
Model (4P)	0.9547	0.0047	164.0	14.8	122.4	276.7

using a dispersion factor [21] of 1.5. Except for a few instances, all observed values were included in the prediction intervals produced by the different models. All three approaches predicted a smaller increase than that observed in the incidence of skin melanoma. The observed number of breast cancer cases was above the upper limit of the prediction interval obtained by models (3) and (3P) showing what would have been expected based on age-specific historical trends. The observed number was within the wider prediction interval produced by model (1). For cancers of the breast, pancreas and lung, model (1) produced a very wide prediction interval; likewise model (4) for cancer of the stomach, indicating a rather limited usefulness of these predictions. The prediction intervals obtained by model (1) were always wider than those for the other models except for skin melanoma.

3.2. Simulation study

The simulation study showed that the distributions of the estimators of the expectation of the predicted value produced by the additive models based on the three approaches are the same and the estimators are unbiased. Appendix B shows in detail that for the same data set, models (1) and (3) always produce the same estimate of the expectation of predicted number of cases. Even though model (3P) does not share this property, the simulated distribution of the expectation for this model (not presented here) was practically the same as for the two other models.

The analysis of simulated distributions of the width of prediction interval for the models highlights the basic differences between the approaches presented here in producing a reliable outcome. Part (a) of Table II illustrates the numerical characteristics of the estimators, and Figure 1 shows graphically the distribution of the width of the prediction interval. Only model (3P) gives prediction intervals, which show a good agreement with the nominal probability of coverage. The other models produce, on average, intervals that are too short. The coverage probabilities of the estimators of models (1) and (3) are relatively high, especially for model (1), due to the large variances of these estimators. Quite appropriately, the estimator of model (3P) is also the most precise among those analysed.

Figure 2 shows the simulated distributions of estimators of the expectation of predicted value for log-linear models. All estimators were biased. The simulated expectation was 429 cases while the means produced by the models (2), (4) and (4P) were 406, 457 and 440, respectively. The

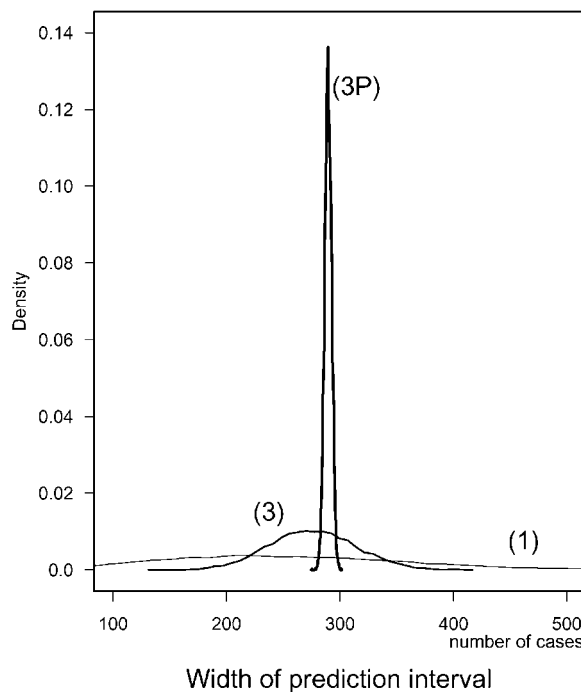


Figure 1. Distribution of the estimator of the width of the prediction interval for the number of breast cancer cases for models (1), (3) and (3P), based on simulations of 5000 observations.

standard error of the estimator for model (4) was the largest and amounted to 42 cases. These quantities for models (2) and (4P) were 29 and 36, respectively. Part (b) of Table II and Figure 3 present the simulated results for estimators of the width of prediction interval for the log-linear models. The estimator for model (2) is clearly biased. Among the two remaining estimators, which have very small coverage error, that for model (4P) is much more precise than that for model (4). On the other hand, the average width of prediction interval for model (4) is smaller than for model (4P). The estimator for model (4) has a disadvantage. It is based on formula (A1) in Appendix A, which involves an exponential transformation, and, as such, is sensitive to changes in the model fit expressed in the formula by the values of $\hat{\sigma}_i^2$. This results in sometimes unrealistic, very wide and useless prediction intervals, even if the model generating the data is perfectly correct. There were 58 such outcomes among 5000 simulations (1.2 per cent) where the width of prediction interval was much bigger than 500 cases, the value never exceeded by any outcome of the other models. Those observations caused the estimator of model (4) to have the biggest standard error among the analysed estimators. This feature disqualifies the model for practical applications and explains why in empirical study for cancer of the stomach (Table I) the width of prediction interval produced by model (4) was a few times larger than that for the other models.

The entire analysis described above was also done for the data tabulated in one-year periods (1954, 1955, ..., 1978) with the period of prediction being 1991, keeping the same structure of age groups. For the data set with five times as many observations, the standard errors of all estimators decreased and the coverage probability of the biased estimators improved. However,

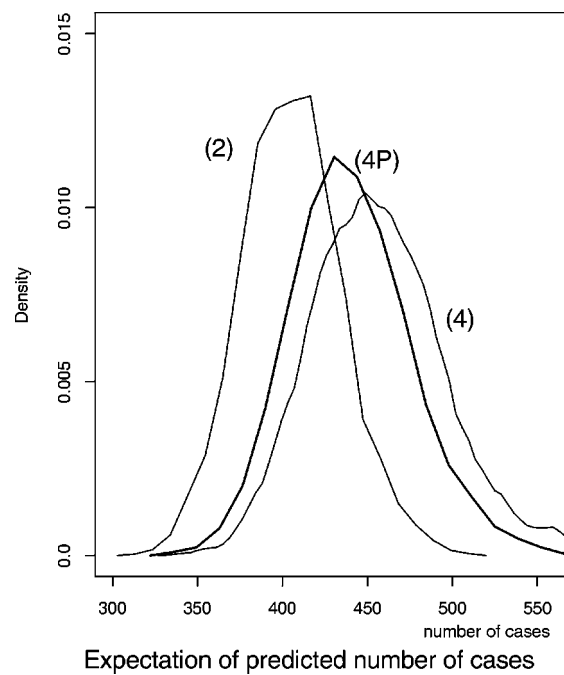


Figure 2. Distribution of the estimator of the expectation of the predicted number of stomach cancer cases for models (2), (4) and (4P), based on simulations of 5000 observations.

the main features presented above – the lowest coverage error and the highest precision of the estimators of models making use of Poisson assumption – did not change.

4. DISCUSSION

Published disease incidence predictions often lack any statement of confidence [22]. Bayesian predictions use credible intervals [23]. The possibility of calculating a prediction interval improves fundamentally the value of the prediction enabling proper conclusions to be drawn from administrative and scientific predictions.

Three methods for performing disease incidence prediction based on a simple interpolation technique were compared in this paper: that based directly on the numbers of observed cases as Poisson variables, and those modelling the age-adjusted and age-specific incidence rates as normally distributed variables. The prediction method based on the Poisson assumption was the most reliable for cancers with both increasing and decreasing incidence patterns. The estimator of the prediction interval for this method had the smallest coverage error and was the most precise and these properties did not depend on the length of the calendar time intervals used for the estimation of the models. Owing to the low precision of the estimators of the width of the prediction interval, the methods making use of a normal approximation can produce unstable, misleading results. For example, the empirical result for breast cancer produced by model (1)

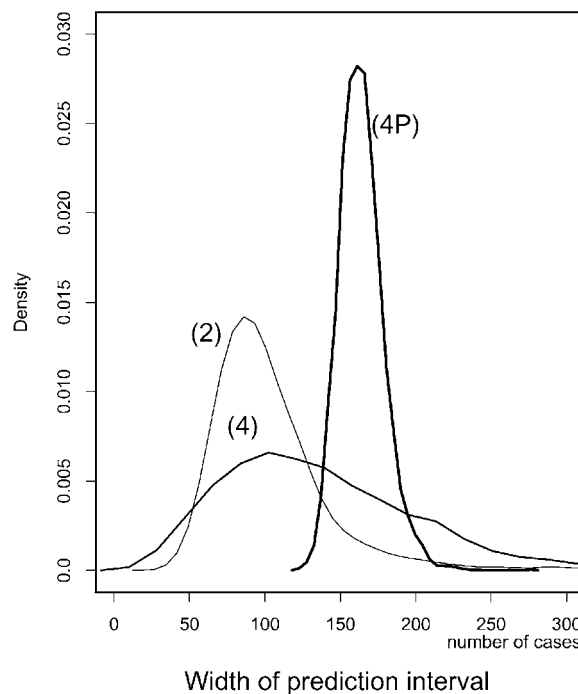


Figure 3. Distribution of the estimator of the width of the prediction interval for the number of stomach cancer cases for models (2), (4) and (4P), based on simulations of 5000 observations.

should be considered as misleading. In 1989–1993, an early detection programme by mammography caused differential age-specific changes in breast cancer incidence in Finland [24], and this was detected only using the age-specific models (3) and (3P). The estimator based on age-adjusted rates is especially questionable when a small number of calendar time intervals is used for the estimation of the models – a small number of degrees of freedom – and improves as this number increases.

Pearson's chi-square statistic [21] is a useful piece of information in the context of prediction. Owing to an increase in the prevalence of tobacco smoking among women in Finland since the 1960s [25], one could expect that the observed number of lung cancer cases would be above the prediction interval, but this was not the case for any approach (Table I). However, the prediction interval based on model (3P) was calculated under overdispersion, making the final prediction interval somewhat wider. The overdispersion is basically a sign of lack of homogeneity in the data, and in our case decreasing the length of the calendar time intervals to one year improved the fit. The fit for the increasingly stratified data did not show a sign of overdispersion and the new predicted average number of cases in 1989–1993 was 217 with a 95 per cent prediction interval 156–280 excluding the average observed number 285 (Table I). A similar exclusion also took place by using model (3) instead of model (3P) based on the annually stratified data.

A good fit does not guarantee a good prediction, but a very good fit should not always lead to a bad prediction, either. However, this happens for models (1)–(5). The general formula for the

variance of prediction of these models [13] (compare also equation (A1) in Appendix A) shows that as the variance of the random component of the model approaches zero, the variance of prediction approaches zero. A perfect fit of these models produces a prediction interval of width zero. The method of calculating the prediction interval [16] for models (3P)–(6P) always guarantees that the variance of the prediction is not smaller than the variance of the future predicted distribution, the natural, intuitively acceptable conclusion.

Among the approaches compared, that taking directly into account the Poisson assumption for number of cases proved to be superior over the methods modelling the age-adjusted or the age-specific incidence rates as normally distributed variables and could thus be recommended for practical applications. It would be advisable to check, however, that the age or regional components are small or sufficiently internally homogeneous to allow the Poisson assumption for the observed numbers of cases. It has now been also demonstrated that asymptotic properties of the estimators of prediction intervals for models (3P) and (4P), derived earlier theoretically, work satisfactorily in practice.

APPENDIX A

Model (4) assumes normal distributions for the stratum-specific error terms, ε_i , implying log-normal distributions for rates, c_{it}/n_{it} . Since those distributions are connected by the logarithmic transformation, an exact relation between them exists [26]. That gives the formula

$$\hat{E}(\hat{c}_T) = \sum_i n_{iT} \exp(\hat{\mu}_i + 0.5\hat{\sigma}_i^2)$$

for the expectation of the predicted number of cases, c_T , where n_{iT} is the age-specific predicted number of person-years for the future period T and where

$$\hat{\mu}_i = \hat{\alpha}_i + \hat{\beta}_i T$$

and

$$\hat{\sigma}_i^2 = \text{var}(\hat{\alpha}_i) + T^2 \text{var}(\hat{\beta}_i) + 2T \text{cov}(\hat{\alpha}_i, \hat{\beta}_i) + \text{var}(\hat{\varepsilon}_i)$$

are, respectively, estimates of the expectation and of the variance of the log-normal distribution for age group i , based on the model. The variance of the prediction is then expressed by the formula

$$\text{var}(\hat{c}_T) = \sum_i n_{iT}^2 \exp(2\hat{\mu}_i + \hat{\sigma}_i^2)(\exp(\hat{\sigma}_i^2) - 1) \quad (\text{A1})$$

Calculating the prediction interval by applying the law of large numbers for summing the age-specific log-normal distributions and using an assumption of normality gives an approximate prediction interval.

APPENDIX B

Let us call the vector of the estimators of the parameters of model (1) $\hat{\beta}^A$. The vector is given as

$$\hat{\beta}^A = [\hat{\alpha}^A, \hat{\beta}^A]' = (X'X)^{-1}X'y \quad (\text{A2})$$

where y is the vector of age-adjusted rates over period t , $t = 1, 2, \dots, N$. Then, for fixed t , $y[t] = \Sigma_i(n_{iT}/n_T)(c_{it}/n_{it})$, $n_T = \Sigma_i n_{iT}$, and the matrix

$$X' = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & 2 & 3 & \dots & N \end{bmatrix}$$

We can then write equation (A2) as

$$\begin{aligned} \hat{\beta}^A &= (X'X)^{-1}X'[\Sigma_i(n_{iT}/n_T)(c_{i1}/n_{i1}), \Sigma_i(n_{iT}/n_T)(c_{i2}/n_{i2}), \dots, \Sigma_i(n_{iT}/n_T)(c_{iN}/n_{iN})]' \\ &= \Sigma_i(n_{iT}/n_T)(X'X)^{-1}X'z_i \end{aligned}$$

where $z_i' = [c_{i1}/n_{i1}, c_{i2}/n_{i2}, \dots, c_{iN}/n_{iN}]$ and because $(X'X)^{-1}X'z_i$, which we here call $\hat{\beta}_i^N$, is the vector $\hat{\beta}^N = [\hat{\beta}_1^N, \hat{\beta}_i^N]'$, of the estimators of the parameters of model (3) for age group i , we can then write

$$\hat{\beta}^A = [\hat{\alpha}^A, \hat{\beta}^A]' = \Sigma_i(n_{iT}/n_T)\hat{\beta}_i^N = \Sigma_i(n_{iT}/n_T)[\hat{\alpha}_i^N, \hat{\beta}_i^N]'$$

It is now seen that the estimated expectation of the predicted number of cases produced by model (1), given by

$$\begin{aligned} \hat{E}(\hat{c}_T) &= (\hat{\alpha}^A + T\hat{\beta}^A)n_T = (\Sigma_i(n_{iT}/n_T)\hat{\alpha}_i^N + T\Sigma_i(n_{iT}/n_T)\hat{\beta}_i^N)n_T \\ &= \Sigma_i n_{iT}(\hat{\alpha}_i^N + T\hat{\beta}_i^N) \end{aligned}$$

is the same as for model (3).

ACKNOWLEDGEMENTS

This work was supported by the Swedish Cancer Foundation, Stockholm Cancer Society, the Ministry of Education in Finland and the Academy of Finland.

REFERENCES

1. Einhorn J. Cancer by the year 2000. Educational requirements for future oncologists. *Acta Oncologica* 1989; **28**:723–728.
2. National Cancer Institute. Cancer control objectives for the Nation: 1985–2000. *NCI Monograph* 1986; **2**:1–105.
3. Schaubel DE, Morrison HI, Desmeules M, Parsons D, Fenton SA. End-stage renal disease projections for Canada to 2005 using Poisson and Markov models. *International Journal of Epidemiology* 1998; **27**:274–281.
4. Prior P, Woodman CBJ, Wilson S, Threlfall AG. Reliability of underlying incidence rates for estimating the effect and efficiency of screening for breast cancer. *Journal of Medical Screening* 1996; **3**:119–122.
5. Coleman MP, Esteve J, Damiecki P, Arslan A, Renard H. *Trends in Cancer Incidence and Mortality*. International Agency for Research on Cancer: Lyon, IARC Scientific Publications No. 121, 1993.
6. Radhakrishna Rao C. *Statistics and Truth*. Council of Scientific and Industrial Research: New Delhi, 1989.
7. Cox DR, Wermuth N. *Multivariate Dependencies-Models, Analysis and Interpretation*. Chapman and Hall: London, 1996.
8. Stephenson RA, Smart CR, Mineau GP, James BC, Janerich DT, Dibble RL. The fall in incidence of prostate carcinoma. *Cancer* 1996; **77**:1342–1348.
9. Hakulinen T, Teppo L, Saxen E. Do the predictions for cancer incidence come true? Experience from Finland. *Cancer* 1986; **57**:2454–2458.
10. Magnus K (ed.). *Trends in Cancer Incidence. Causes and Practical Implications*. Hemisphere: Washington, 1982.
11. Breslow N, Day N. *Statistical Methods in Cancer Research Vol. II*. International Agency for Research on Cancer: Lyon, IARC Scientific Publications No. 82, 1987.
12. Boyle P, Parkin DM. Statistical methods for registries. In *Cancer Registration: Principles and Methods*, Jensen OM, Parkin DM, MacLennan R, Muir CS, Skeet RG (eds.) International Agency for Research on Cancer: Lyon, IARC Scientific Publications No. 91, 1991.

13. Fleiss L. *Statistical Methods for Rates and Proportions*. Wiley: New York, 1973.
14. Weisberg S. *Applied Linear Regression*. Wiley: New York, 1985.
15. Judge GG, Griffiths WE, Hill RC, Lutkepohl H, Lee T. *The Theory and Practice of Econometrics*. Wiley: New York, 1985.
16. Hakulinen T, Dyba T. Precision of incidence predictions based on Poisson distributed observations. *Statistics in Medicine* 1994; **13**:1513–1523.
17. Dyba T, Hakulinen T, Päivärinta L. A simple non-linear model in incidence prediction. *Statistics in Medicine* 1997; **16**:2297–2309.
18. Clayton D, Schifflers E. Models for temporal variation in cancer rates II: Age-period-cohort models. *Statistics in Medicine* 1987; **6**:469–481.
19. Dyba T. Confidence and prediction intervals for disease incidence using Glim. *Glim Newsletter* 1995; **25**(25):27–32.
20. Hakulinen T, Kenward M, Luostarinen T, Oksanen H, Pukkala E, Söderman B, Teppo L. *Cancer in Finland in 1954–2008*. Cancer Society of Finland: Helsinki, 1989, publication no. 42.
21. McCullagh P, Nelder JA. *Generalized Linear Models*. Chapman and Hall: London, 1989.
22. Hakulinen T. Cancer projection methods in Europe. In *Proceedings of the Canadian Cancer Projection Workshop June 8–10, 1991*, MacLauhglin J, Morgan P, Mao Y. (eds.) Health and Welfare Canada: Toronto, 1992; 1–14.
23. Berzuini C, Clayton D. Bayesian analysis of survival on multiple time scales. *Statistics in Medicine* 1994; **13**:823–838.
24. Hakulinen T. The future cancer burden as a study subject. *Acta Oncologica* 1996; **35**:665–670.
25. Saxen E, Teppo L, Hakulinen T. Epidemiology of lung cancer in Scandinavia. In *Air Pollution and Cancer in Man*. Mohr U, Schmähl D, Tomatis L (eds). International Agency for Research on Cancer: Lyon, IARC Scientific Publications No. 16, 1977; 217–228.
26. Atchison J, Brown JAC. *The Lognormal Distribution*. Cambridge University Press: Cambridge, 1969.