

# Analysis of Breast Cancer Survival in Local Health Authorities

Joseph Lynch<sup>1</sup> and Gilbert MacKenzie<sup>1</sup>

<sup>1</sup> Centre of Biostatistics, University of Limerick, Ireland

**Abstract:** Kaplan-Meier analysis of a large breast cancer dataset is carried out under all-cause and cause-specific survival. The results are compared with a variety of model-based analyses, including Cox's Proportional Hazard (PH) model and its Gamma frailty variant, along with the non-PH Generalised Time-Dependent Logistic Model (GTDL) and its Gamma frailty variant.

**Keywords:** Kaplan-Meier; cause-specific; all-cause; PH/non-PH; Gamma frailty

## 1 Introduction

Coleman (1999) reported that North Staffordshire Local Health Authority (LHA) was ranked last of 99 LHAs in England & Wales with respect to breast cancer survival. His report was based on a relative survival approach which did not take account of *case-mix* factors and he analysed incident cases diagnosed between 1991-1993. We re-analyse an augmented dataset from the West Midlands of England, including North Staffordshire, by more traditional methods and report on the resulting *case-mix* adjusted league table.

The population data analysed comprise 15,516 incident cases of cancer of the female breast diagnosed in the West Midlands, UK, between 1991-1995 and followed-up to the end of 2001. Survival time is defined as the time in years from diagnosis to death or censoring. Both cause-specific and all-cause definitions are used. The cause-specific definition refers to deaths in which breast cancer is registered as the primary cause - other outcomes being regarded as censored (at their time of occurrence). In contrast the all-cause refers to death from all causes including breast cancer.

## 2 Model Definitions

We consider several models of increasing complexity; the proportional hazard (PH) model of Cox (1972), the Gamma frailty variant (Hougaard, 1994), the Generalised Time-Dependent Logistic Model (MacKenzie 1996, 1997) and the Gamma frailty variant discussed by Blagojevic, MacKenzie

& Ha (2003). The models are defined in order by:

$$\begin{aligned}\lambda(t|x) &= \lambda_0(t)\exp(x\beta) \\ \lambda(t|u, x) &= u\lambda_0(t)\exp(x\beta) \\ \lambda(t|x) &= \lambda p(t|x) \\ \lambda(t|u, x) &= u\lambda p(t|x)\end{aligned}\tag{1}$$

where  $\lambda_0(\cdot)$  is an unspecified baseline hazard function,  $\beta$  is a  $p \times 1$  vector of regression parameters associated with fixed covariates,  $x$  and  $U \sim \text{Gamma}$  with  $E(U) = 1$  and  $V(U) = \sigma^2$ . Because some of the covariates studied do not obey the PH assumption we also adopted the non-PH GTDL model with  $p(t_i|x_i) = \exp(t_i\alpha + x_i'\beta)/1 + \exp(t_i\alpha + x_i'\beta)$  for  $i = 1, \dots, n$  subjects.

### 3 Preliminary Results

Below we report the major findings from the PH model with which the Epidemiologists involved in the study are most familiar. Further results from the more sophisticated models will be presented in the main paper.

The data were obtained from the West Midlands Cancer Intelligence Unit (the Cancer registry) in Birmingham, UK. There are 10 major covariates of interest including: age, diagnosis basis, stage, grade, morphology, whether or not the cancer was screen-detected, Townsend score (a measure of deprivation), year of diagnosis, Local Health Authority (14 including one unknown category) and treatment. All covariates were categorical.

Figure 1 shows the overall KM survival rates, the upper curve refers to cause-specific survival (mean = 8.33 years) and the lower curve to all-cause survival (mean = 7.43 years), so that (typically) all-cause mortality is higher, leading to shorter survival.

From the KM analysis, using LHA as a factor, we were able, somewhat surprisingly, to produce a result analogous to Coleman's in this extended data set, which had more cases recruited and a longer follow-up period. North Staffordshire was again bottom of the league - this time the West Midland's league. As with Coleman's analysis, no covariates were used at this stage.

In the Cox analysis, all ten of the covariates were statistically significant. The major factors influencing survival were, in order, stage, treatment, grade, age and tumour morphology. These highly significant results were not unexpected, given the very large sample size involved, and it is arguable that the magnitudes of the adjusted odds-ratios are perhaps a more interesting measure of the joint impact of the ten covariates.

We addressed the underlying variable selection issue by generating 100 bootstrap samples from the original data and re-fitting the model using a variety of stepwise algorithms and cut-off points including different values of  $t = \hat{\beta}/\text{se}(\hat{\beta})$  and various Odds Ratios. The stability of variable selection

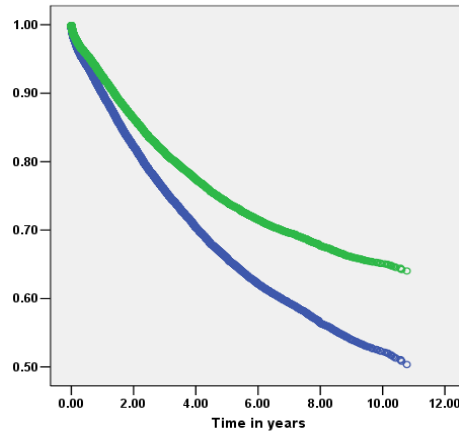


FIGURE 1. KM Breast Cancer Survival: cause specific=upper curve, all cause=lower curve

was re-assuringly good in the cause-specific analysis. For example, when  $|t| > 2$  all variables were selected in 100% of the bootstrap samples and when  $|t| > 3$ , Health Authority was selected in 78% of samples and social class in 94%. All other covariates were selected in 100% of the samples. In the all-cause analysis, only Health Authority exhibited variation being selected in 96% of the samples. Accordingly, for the purposes of this paper, we elected to adjust the LHA rates for the remaining 9 covariates.

Table 1 shows the resulting 5-year survival league tables for all-cause (unadjusted and adjusted,) with approximate standard errors. There is wide variation between the unadjusted and adjusted ranking results. However these findings should be interpreted cautiously as: (a) the quantitative differences are small and (b) when North Staffordshire is regarded as the reference category in the regression analysis only Coventry and Birmingham have significantly better survival.

Thus we see that Coleman's claim - that North Staffordshire is at the bottom of the League is not quite justified when key *case-mix* factors are taken into account. However, clearly, breast cancer survival in North Staffordshire was not optimal over the study period.

## 4 Discussion

This is a large study and we have only just begun to scratch the surface. The preliminary findings suggest that covariate adjustment is required for valid interpretation. However, here, the PH model is merely an approximation to the truth since in the course of the analysis we found that several important covariates did not obey the PH assumption. In the main paper we shall

TABLE 1. All-cause 5-year Survival League tables

LHA	$\hat{S}_{KM}(t = 5)$	95%CI	LHA	$\hat{S}_{PH}^*(t = 5 \bar{x})$
Solihull	0.71	(0.66, 0.73)	Unknown	0.783
Worcester	0.68	(0.65, 0.69)	Birmingham	0.763
Hereford	0.68	(0.63, 0.71)	Hereford	0.753
Shropshire	0.68	(0.65, 0.69)	Coventry	0.752
Warwick	0.68	(0.65, 0.69)	Warwick	0.732
Wolverhampton	0.67	(0.64, 0.72)	Sandwell	0.727
Coventry	0.67	(0.62, 0.68)	Shropshire	0.727
South Staffs	0.65	(0.62, 0.67)	Solihull	0.727
Walsall	0.65	(0.60, 0.68)	Wolverhampton	0.722
Unknown	0.65	(0.57, 0.96)	Worcester	0.716
Birmingham	0.65	(0.62, 0.68)	Dudley	0.706
Dudley	0.65	(0.63, 0.68)	Walsall	0.706
Sandwell	0.63	(0.60, 0.69)	North Staffs	0.696
NStaffs	0.58	(0.57, 0.63)	South Staffs	0.686

Confidence Intervals for adjusted survival models are pending \* $\bar{x}$  is the West Midlands mean.

compare the findings obtained by fitting the non-PH models, discuss the value of frailty in this context and comment on the process of obtaining adjusted survival curves. These analyses are now in hand.

### References

- Coleman, *et al* (1999) Cancer Survival in the Health Authorities of England up to 1998; A report prepared for the National Health Services Executive .
- Cox DR (1972) Regression models and life-tables (with discussion). *J. R. Statist. Soc. B*, **34**, 187-220.
- Hougaard, P. (1994). Heterogeneity Models of Disease Susceptibility, with Applications to Diabetic Nephropathy. populations. *Biometrics*, 50, 1178-1188.
- MacKenzie, G. (1996) Regression models for survival data: the generalised time dependent logistic family. *JRSS Series D*, 45, 21-34.
- MacKenzie, G. (1997) On a non-proportional hazards regression model for repeated medical random counts. *Statistics in Medicine*, 16, 1831-1843.
- Blagojevic M., MacKenzie G. and Ha I.D. (2003) A Comparison of non-PH & PH - Gamma frailty models. *IWSM 2003,pp 39-43*.