

# THE UNIVERSITY *of York*

## *Discussion Papers in Economics*

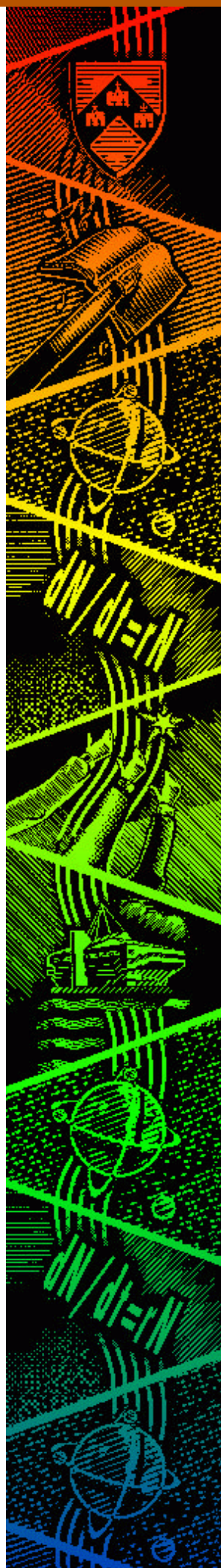
No. 2000/51

The Role of Tobacco Taxes in Starting and Quitting Smoking:  
Duration Analysis of British Data

by

Martin Foster and Andrew M Jones

Department of Economics and Related Studies  
University of York  
Heslington  
York, YO10 5DD



# The role of tobacco taxes in starting and quitting smoking: duration analysis of British data

Martin Forster and Andrew M. Jones

*Department of Economics and Related Studies, University of York, Heslington, York YO10 5DD, UK.*

*e-mail: mf8@york.ac.uk, amj1@york.ac.uk*

**Summary.** The annual five per cent in real terms increase in tobacco taxes proposed in the recent White Paper on smoking has reaffirmed the commitment of successive United Kingdom governments to above inflation increases in tobacco taxation to encourage people to stop smoking. This paper presents evidence on the determinants of starting and quitting smoking using data from the British *Health and Lifestyle Survey* and is the first to identify tax elasticities for starting and quitting smoking using British data. Self-reported individual smoking histories are coupled with a long time series for the tax rate on cigarettes to construct a longitudinal data set. Estimates are obtained for the impact of above inflation tax rises on the age of starting smoking and the number of years of smoking. The estimates of the tax elasticity of the age of starting smoking are +0.16 for men and +0.08 for women. The estimates of the tax elasticity of quitting are -0.60 for men and -0.46 for women. These are robust to different specifications.

*Keywords:* Smoking initiation and cessation; Tobacco taxes; Duration analysis

## 1. Introduction

This paper presents new evidence on the determinants of starting and quitting smoking using duration data from the British *Health and Lifestyle Survey* (HALS). Self-reported data on individual smoking histories coupled with the availability of a long time series for the tax on cigarettes allow us to construct a longitudinal data set in which the tax on cigarettes is treated as a time varying covariate. This overcomes the problem of the lack of cross section variation in prices that has plagued previous studies of smoking in Britain. Only a handful of previous studies have used duration analysis to model the hazard rates of starting and quitting and all of these have used data from the United States. This is the first study to identify tax elasticities for starting and quitting in Britain, where the tax elasticities measure the impact of a proportionate change in the real tax on cigarettes on the age of starting smoking and the number of years an individual smokes for.

Results for the age of starting smoking are reported for a split population duration model like that used by Douglas and Hariharan (1994). Results for the number of years smoked prior to quitting are reported for Cox, Weibull and generalized gamma models. All of the models are estimated separately for men and women, with extensive diagnostic tests used to guide model specification. Results and specifications are compared to those in the U.S. literature. We find that, for the starting models, the split population log-logistic duration/probit participation specification is well specified for men but not for women. For the quitting models, the generalized gamma specification is preferred to the specifications used by Douglas (1998) and Tauras and Chaloupka (1999), although tax elasticities of quitting show little variation across models.

To assess the robustness of the estimated tax elasticities for starting and quitting we carry out sensitivity analyses. For the models of starting we assess the effect of recall bias with respect to the age at which an individual started smoking by rescaling the duration variable to measure calendar time and investigating systematic reporting bias implied by the hazard function. We find little evidence of recall bias using this method. We also compare our results from the split population models with ‘two-part’ specifications with a probit model for starting and a log-logistic model for age of starting estimated on the sub-sample of ‘starters’. Estimates of the tax elasticity in these log-logistic models are very similar to the split population models. For the models of quitting we compare different semiparametric and parametric specifications of the baseline hazard function, compare discrete and continuous time specifications, and allow for the influence of unobserved heterogeneity using a mixture model. We also assess recall bias by rescaling the duration variable for quitting to investigate the quitting hazard function by calendar year. We find strong evidence of a ‘heaping’ effect, whereby responses are clustered around five and ten year points prior to the survey date. Our models are adapted in response to this evidence using the methods of Torelli and Trivellato (1993). For comparison with the U.S. literature, we also present estimates that include measures of past smoking. The sensitivity analyses suggest that the estimated tax elasticities of quitting are robust to all of these factors.

In our benchmark models, the estimated tax elasticity of the age of starting is +0.16 for men and +0.08 for women. The estimated tax elasticity of the number of years of smoking before quitting is -0.60 for men and -0.46 for women. Parental smoking increases the probability of becoming a smoker and reduces the age of starting but has no significant effect on quitting. Those with higher educational attainment are less likely to start and start later. There is a clear socio-economic gradient in success in quitting, and those with higher educational attainment smoke for shorter durations.

Our estimated tax elasticities relate to the impact of above inflation tax rises on the number of years smoked by current smokers and our results are relevant for current U.K. government public health initiatives. Recent White Papers on smoking and on public health (Department of Health, 1998b, 1999) have stressed the central role of tobacco control in the policy goals of improving population health and of reducing socio-economic inequalities in health. Since the early 1990s successive governments have had a commitment to annual increases in the real level of tobacco taxes, to achieve health policy objectives and encourage people to stop smoking. In the Budget of July 1997, the Chancellor of the Exchequer announced that in future budgets tobacco duties would increase on average by at least 5 per cent in real terms, and the most recent White Paper on smoking (Department of Health, 1998b) reaffirmed this commitment. More recently a commitment has been made to hypothecate the revenues from the duty escalator to fund the National Health Service.

## **2. Previous evidence on starting and quitting**

Grossman (1999) provides a comprehensive review of the evidence of the influence of prices on substance use and abuse, and recent surveys of the literature on starting and quitting smoking are provided by Douglas (1998) and Tauras and Chaloupka (1999). There are two main approaches to modelling smoking behaviour: those that treat starting and quitting as binary events within a discrete choice framework, and those that use duration models.

Discrete choice models of smoking among teenagers typically report greater price responsiveness than among the population at large (see, for example, Lewit et al. (1981),

Lewit and Coate (1982), Chaloupka and Grossman (1996), Chaloupka and Wechsler (1997), Evans and Farrelly (1998) and Harris and Chan (1999)). But these often use the level of consumption or the overall participation rate as their dependent variable, rather than the age of starting smoking. A number of retrospective studies have examined the effects of demographic variables, health and past smoking, on the propensity to start and quit smoking. These include Jones (1989, 1994), Sander (1995), Yen and Jones (1996), Hsieh (1998) and Dorsett (1999).

An alternative to estimating discrete choice models is to use duration analysis, as in Douglas and Hariharan (1994), Douglas (1998) and Tauras and Chaloupka (1999), all of whom use U.S. data. Douglas and Hariharan (1994) use a split population model like that used by Schmidt and Witte (1989) to model the hazard of starting smoking, using data from the 1978 and 1979 Smoking Supplements to the U.S. National Health Interview Survey. They find evidence that those with higher lifetime educational attainment are less likely to start smoking and, if they do start, start later. Women are less likely to start and do so at later ages. No evidence is found that higher real prices of cigarettes (measured at the time the individual was aged 18 and as the change in the real price between the age of 15 and 18) significantly reduce the probability of starting smoking or significantly increase the age at which individuals start.

Douglas (1998) considers the hazards of both starting and quitting using the 1987 U.S. National Health Interview Survey. He uses a split population model based on an ordered probit, which distinguishes between those who never start smoking, those who start and will eventually quit and those who start and never quit. The price of cigarettes is treated as a time varying covariate. He uses a log-logistic function for the starting hazard and a Weibull model for the quitting hazard. The price of cigarettes during the teenage years has no significant effect on the hazard of starting smoking or the probability of becoming a smoker and the number of years an individual smokes is found to have an approximately unitary elasticity with respect to the price of cigarettes.

Tauras and Chaloupka (1999) estimate Cox proportional hazard models of quitting using longitudinal data from the U.S. Monitoring the Future Surveys consisting of a representative sample of high school seniors split by gender. They estimate average price elasticities of the quitting hazard to be significant at +1.12 and +1.19 for young adult men and women respectively.

### **2.1. Comment**

Tauras and Chaloupka (1999) express concern about using retrospective data sets in analyses of smoking behaviour. They argue that recall bias can be a serious problem in such studies, leading to errors in individuals' reported age of starting and quitting smoking which can bias parameter estimates. Further bias may occur if an individual's current state of residence is used to impute past prices of cigarettes - cross-state variation in cigarette prices occurs in the U.S. and data for individuals who move states will be subject to measurement error. Shmueli (1996) has criticised the use of measures of self-assessed health in retrospective analyses of quitting behaviour. This is because of problems with unobservable heterogeneity bias, which may also be a problem if measures such as previous peak consumption are used to predict quitting behaviour.

An additional problem with previous studies - both retrospective and prospective - is their lack of diagnostic tests to assess the fit of the models. When fitting their Cox proportional hazards models, Tauras and Chaloupka (1999) do not report whether the

assumption of proportional hazards is valid for the data and parameters of their model. If the assumption is not valid, parameter estimates of elasticities will be inconsistent. Douglas and Hariharan (1994) provide a graphical assessment of their split population model but Douglas (1998) does not report the specification adequacy of his model.

This paper addresses these issues as follows. Because U.K. tobacco taxes are set at a national level, we do not encounter the problems of state-level variation that occur with U.S. data. To avoid the problems highlighted by Shmueli (1996), we estimate parsimonious models that include only independent variables that are likely to be exogenous to the individual such as gender, ethnic origin, parental smoking, education level and social class. We check the sensitivity of the parameter estimates to variation in the tax data used in the starting and quitting equations and the inclusion of measures of previous smoking behaviour. We test the models using a number of diagnostic tests from the econometric and biostatistics literature and we check the sensitivity of the quitting equations using various parametric and semi-parametric forms of the hazard function. Recall bias is tested using a rescaling of the duration variable by calendar year in a way previously used by Tunali and Pritchett (1997), which allows us to check for systematic bias in reporting by calendar year. We then adjust our models of quitting to take account of the evident ‘heaping’ effect using the methods considered by Torelli and Trivellato (1993).

Ideally this kind of analysis should use a prospective longitudinal data set, as advised by Tauras and Chaloupka (1999). However, such high quality data is not available for the analysis of smoking behaviour over a long time horizon. By using retrospective data we can begin our period of analysis in 1920 and take account of long-run variations in cigarette taxes. This allows us to assess the effect of both tax effects and the separate influence of the health scares associated with smoking that occurred during the 1960s. Our diagnostic checks and sensitivity analyses are intended to yield robust estimates of the impact of the tax on cigarettes, as well as the effect of demographics and parental smoking behaviour, on the decisions to start and quit smoking.

### 3. Data and sampling

*The Health and Lifestyle Survey* (HALS) is a study designed to record the lifestyles, personal circumstances and the physical and mental health of a large representative sample of individuals aged 18 and over living in households in England, Scotland and Wales in 1984 (Cox et al., 1987). The focus in this paper is on the first wave of the survey and the smoking histories of participants.

HALS consisted of two home visits. In the first, an interviewer questioned the participant about self-reported health, health attitudes and health-related lifestyle such as diet, exercise, smoking and alcohol consumption. Data were also collected on demographic characteristics, employment status, qualifications and household income. In the second visit, a nurse took various measures of physiological and cognitive function and gave participants a self-completion questionnaire to collect information on personality and mental health.

The HALS sample who participated in the original home interview survey numbered 9003 individuals, a response rate of 73.5% of those initially randomly selected. Cox et al. (1987) argue that the study sample is ‘a good and representative sample of the population’. The analysis in this paper uses self-reported information on individuals’ smoking histories to construct duration variables, and is therefore retrospective. The following sections describe the calculation of the duration data, the tax data and the other covariates used in the

estimation.

### 3.1. Duration data

The *Health and Lifestyle Survey* contains retrospective information on whether or not an individual started smoking. It also provides information that separates non-smokers into those who have never smoked and those who class themselves as ex-smokers, allowing the analysis to be extended to distinguish between starting and quitting. These discrete choices are augmented by data on the age of starting smoking and the number of years an individual has smoked for. These can be interpreted as ‘failure times’ and estimated by duration analysis. In contrast to Douglas’s (1998) recent paper on starting and quitting smoking, we model the starting and quitting processes separately.

In our data a smoker is defined as someone who has smoked at least one cigarette per day for a minimum of six months. For the analysis of starting, the self-reported variables FAGAGE and EXFAGAGE are used to measure the age of starting in years for those individuals who have smoked at some point in their lives.† The indicator variable  $c = 1$  denotes a current or ex-smoker at HALS,  $c = 0$  otherwise. This is a self-reported measure.‡ Analysis of the duration of smoking is carried out on the sub-sample of individuals who had started smoking ( $c = 1$ ). The variable  $\text{SMOKE\_YEARS}_i$  is calculated for individual  $i$  as:

$$\text{SMOKE\_YEARS}_i = \text{INTERVIEW\_DATE}_i - \text{START\_DATE}_i - \delta_i \text{QUIT}_i.$$

$\text{INTERVIEW\_DATE}_i - \text{START\_DATE}_i$  measures the number of days between the date of the interview§ and the date individual  $i$  started smoking.  $\text{QUIT}_i$  measures the number of days since the individual quit for ex-smokers at the time of HALS (for whom  $\delta_i = 1$ ).¶ Individuals who have quit smoking represent completed spells and those who are still smoking at the time of HALS represent censored spells. This duration variable is then rounded to the nearest year.

Using these measures of duration, we can also identify the calendar year in which an individual started smoking and the calendar year in which they quit smoking. These can be linked to time series data on tax rates. For the HALS sample the year of starting ranges from 1909 to 1984 and for quitting it ranges from 1913 to 1985. Matching this information with annual tax data provides scope for exploiting a long time series with sufficient variability in the tax on cigarettes to identify tax elasticities of starting and quitting.

### 3.2. Tax data

Empirical work on smoking in the United States has been able to exploit state and local variations in tobacco taxes to identify cross section variation in the price of cigarettes (see

†These variables relate to questions 58(a) and 60(a) in the survey questionnaire which ask ‘How old were you when you started to smoke cigarettes?’ to current and ex-smokers.

‡This is constructed from questions: 55(a) - ‘Now, do you regularly smoke cigarettes, that is, do you regularly smoke at least one cigarette a day?’ and 55(c) - ‘Have you ever smoked at least one cigarette a day for as long as six months?’.  $c = 1$  for anyone who responds ‘yes’ to either question.

§The date of interview for each respondent is not provided in the HALS data set, but Cox (1999) provided us with information on estimated dates calculated from other records in the survey. All variables are converted to a common time scale using the ‘elapsed time’ construction in STATA release 6.

¶ $\text{QUIT}_i$  is computed from the survey variable EXFAGAN which relates to question 60(f) - ‘How long ago did you completely stop smoking cigarettes?’.

Chaloupka (1991), Douglas and Hariharan (1994) and Douglas (1998)). In Britain, the tax rate for cigarettes is set each year by central government and the real price of cigarettes does not vary across regions to the degree that it does in the U.S.. This has prevented comparable analyses of the tax/price elasticities of starting and quitting in Britain. In this study we solve the problem by using time series variation coupled with retrospective data on smoking durations.

We use the ‘tax per cigarette’, calculated using the total receipts from tobacco duty as a share of total sales volume as a proxy for the real price of cigarettes. The second edition of *U.K. Smoking Statistics* (Wald and Nicolaides-Bouman, 1991) provides an unbroken annual series for the financial years 1920-21 to 1989-90 for the total receipts from tobacco duty for the UK.|| Wald and Nicolaides-Bouman also present a series for total annual sales of manufactured cigarettes (numbers in millions) for the UK for the calendar years 1905 to 1987.\*\* Tax receipts are divided by sales volume to give an annual series for the tax per cigarette in constant 1913-14 prices. These data are shown in Fig. 1. The data on tax receipts in Wald and Nicolaides-Bouman also allow us to construct annual data for 1905-1913, but suitable data are not reported for the years 1914-1919. Only 209 individuals from our full sample of 9003 started smoking before 1920 and we exclude them from our main analysis. However, as part of the sensitivity analysis, models of the hazard of starting are re-estimated with the full tax series for 1905-1985. Values for 1914-1919 are set equal to the average of the 1913 and 1920 values. Only two observations are lost from the analysis of the hazard of quitting by omitting the tax data for the years 1905-1919.

As tobacco duty has been a high proportion of the price of cigarettes throughout the century this may be a reasonable proxy for the real price of cigarettes. However the measure will be contaminated by variations in the share of tax in the full price of cigarettes over time.†† It is possible to construct an RPI series for cigarettes for the period 1948-1985. However, for 1948-1951 this aggregates alcohol and tobacco, for 1952-1974 it covers all tobacco and from 1974 onwards it covers cigarettes. The correlation coefficient between the natural logarithm of the tax per cigarette and the natural logarithm of the price of cigarettes using the RPI data is 0.59.

The tax on cigarettes is the relevant policy instrument for the government if it wishes to use fiscal policy to influence smoking. In this sense, our estimates can be viewed as elasticities for the ‘policy response’, relating changes in rates of taxation to their effect on starting and quitting. To analyse the impact of the tax on cigarettes on starting and quitting smoking these data need to be mapped to the individual observations in the HALS data.

||For 1920-21 to 1975-76 this is measured in 1913-14 prices and for 1976-77 to 1989-90 it is measured in 1974 prices. The constant price series are based on the all items Retail Price Index (RPI), full details of which are given in Wald and Nicolaides-Bouman (1991). Our results may be sensitive to the choice of this deflator.

\*\*Calendar year is matched with the first year of the fiscal year. That is, the 1920 volume data is matched with the 1920-21 tax data. This is intended to give the greatest possible overlap (nine months) between the series.

††In a study of smoking among U.S. College students, Chaloupka and Wechsler (1997) use information on state excise tax rates as an alternative measure of cigarette prices as a response to the potential endogeneity of the state level price data. Encouragingly they report that ‘estimated elasticities based on the models using tax data as an instrument for price are generally very similar to those based on models using the price itself’. As our models include a flexible time trend this should pick up any secular trend in the relationship between the tax and price per cigarette.



**Fig. 1.** Index for tax per cigarette.

### 3.3. Other covariates

Because of the potential problems associated with predicting past behaviour as a function of individual characteristics that are measured at the time of the HALS survey we use a parsimonious set of exogenous covariates for the starting and quitting models. This attempts, as far as possible, to use covariates that were exogenously determined prior to an individual's starting or quitting decision and therefore avoids covariates, such as health and past smoking status, that may be prone to unobservable heterogeneity bias. Summary statistics for the variables used in the models are presented in Table 1.

Educational status is measured by the highest qualification attained by the individual, running from the lowest through to the highest.†

†We follow the classification of qualifications used in the HALS questionnaire to allow for the cross-time equivalence of levels of educational attainment. Our classifications are as follows:

- NO\_QUALIFICATION: no qualifications;
- CSE/O\_LEVEL: CSE grades 1-5, GCE 'O' level, School Certificate, Scottish Lower, City and Guilds Craft/Ordinary level;
- A\_LEVEL: GCE 'A' and 'S' level, Higher Certificate, Matriculation, Scottish Higher;
- HND: ONC/OND/City and Guilds Advanced/Final level, HNC/HND/City and Guilds Full Technological Certificate, RSA/other clerical and commercial;
- DEGREE: degree, including higher degree;
- OTHER: teacher training, nursing, professional and other work-related qualifications.



**Table 1.** Variable definitions and sample means

variable		full sample		starting		quitting	
		men	women	men	women	men	women
AGESTART	age of starting smoking	-	-	-	-	15.45	18.54
START	=1 if started smoking	-	-	0.66	0.52	1.00	1.00
<i>Social class</i>							
SOCIAL_CLASS_1s	social class 1/student	0.06	0.06	-	-	0.03	0.04
SOCIAL_CLASS_2	social class 2	0.22	0.23	-	-	0.20	0.20
SOCIAL_CLASS_3a	social class 3 non-manual	0.37	0.34	-	-	0.40	0.38
SOCIAL_CLASS_3	social class 3 manual			omitted regressor			
SOCIAL_CLASS_4	social class 4	0.17	0.17	-	-	0.20	0.19
SOCIAL_CLASS_5n	social class 5	0.07	0.06	-	-	0.07	0.07
<i>Highest qualification</i>							
NO_QUALIFICATION	no qualification	0.46	0.52	0.45	0.51	0.63	0.67
CSE/O_LEVEL	cse/O level			omitted regressor			
A_LEVEL	A level	0.05	0.05	0.05	0.05	0.03	0.04
HND	higher national diploma	0.12	0.10	0.12	0.10	0.08	0.07
DEGREE	degree	0.13	0.12	0.13	0.12	0.09	0.08
OTHER	other qualifications	0.06	0.03	0.07	0.03	0.08	0.03
<i>Ethnic origin</i>							
WHITE	white/European			omitted regressor			
ASIAN	Indian/Pakistani/Bangladeshi	0.02	0.01	0.02	0.14	0.01	0.00
AFRO-CARIBBEAN	black/African/West Indian	0.01	0.01	0.01	0.01	0.01	0.00
OTHER_NON-WHITE	other non-white	0.01	0.01	0.01	0.01	0.00	0.01
<i>Parental smoking</i>							
PARENTS_NON-SMOKERS	neither parent smoked			omitted regressor			
MOTHER_SMOKED	only mother smoked	0.06	0.07	0.06	0.07	0.05	0.06
FATHER_SMOKED	only father smoked	0.48	0.45	0.48	0.46	0.58	0.50
PARENTS_SMOKED	both parents smoked	0.32	0.31	0.33	0.33	0.26	0.32
YEAR	years since 1920	-	-	35.38	36.54	41.24	45.34
Number of observations		3905	5098	3737	4861	2480	2482
Number of failures				2460	2508	1176	938

Additional covariates chosen to control for other potential influences on smoking are social class (included in the quitting models only), ethnic origin, gender, and parental smoking. The fact that all the variation in taxes is attributable to variation across calendar years raises an identification problem for separating the time trend and tax effects. Our solution is to use a 4th order polynomial to impose a smooth but flexible time trend and to identify tax effects by variations around this trend. We therefore created a variable YEAR to measure the number of years since 1920 and included a quartic polynomial in YEAR to capture any trends in the data independent of the tax effects.‡ The time trend is intended to capture the secular trend in smoking participation, particularly associated with the cumulative impact of increased awareness of the health risks of smoking since the early 1960s and changing perceptions of smoking over the period.

## 4. Methods

### 4.1. Starting

For the smokers in the sample we use the reported age of starting and the duration data can be interpreted as a complete spell. However the sample also contains individuals who have not started. In a parametric duration model these observations are interpreted as incomplete spells, and it is assumed that all of these individuals will eventually ‘fail’ and start smoking. They are classed as ‘right censored’ at the time of the survey.

In their analysis of U.S. data on the age of starting smoking, Douglas and Hariharan (1994) argue that standard duration analysis techniques may not be appropriate and that a split population model should be used. The theory and logic behind the split population duration model used to analyse starting is explained in detail in Schmidt and Witte (1989) who apply it to the study of criminal recidivism, and Douglas and Hariharan (1994) and Douglas (1998), who apply it to the study of starting smoking. In the split population duration model, the duration process applies only to those individuals who are predicted eventually to start smoking. Defining  $s = 1$  for an individual who will eventually start smoking and modelling eventual failure using a probit specification, yields:

$$\begin{aligned} P(\text{eventually start smoking}) &= P(s = 1) = \Phi(\boldsymbol{\alpha}'\mathbf{z}_i), \\ P(\text{never start smoking}) &= P(s = 0) = 1 - \Phi(\boldsymbol{\alpha}'\mathbf{z}_i), \end{aligned}$$

where  $\mathbf{z}_i$  is a vector of time invariant covariates,  $\Phi$  is the cumulative density function for the standard normal distribution and  $\boldsymbol{\alpha}$  is a parameter vector. The probability of starting smoking at a given time  $t$  is then defined conditional upon eventually starting.

Following Douglas and Hariharan (1994), and the evidence of a plot of the product-limit estimate of the hazard function for the age of starting, we choose a log-logistic distribution to model duration. The probability density function,  $f(\cdot)$ , and the survival function,  $S(\cdot)$ , of the log-logistic distribution for those individuals who eventually start smoking are, respectively, (Greene, 1993):

$$f(t|s = 1; \mathbf{x}_i(t)) = \frac{\lambda^{\frac{1}{\gamma}} t^{\frac{1}{\gamma}-1}}{\gamma \left[1 + (\lambda t)^{\frac{1}{\gamma}}\right]^2},$$

‡Section 5.2.1 shows how this approach is adapted to allow for the possibility of recall bias in the durations. The sensitivity analysis experiments with different specifications of the time trend.

$$S(t|s = 1; \mathbf{x}_i(t)) = \frac{1}{1 + (\lambda t)^{\frac{1}{\gamma}}},$$

where  $\lambda = \exp(-\boldsymbol{\beta}'\mathbf{x}_i(t))$ ,  $\mathbf{x}_i(t)$  is a vector of time variant and time invariant covariates and  $\gamma$  is a scale parameter.

The contribution to the log-likelihood function for the split population model then becomes, for individual  $i$ :

$$c_i \ln [\Phi(\boldsymbol{\alpha}'\mathbf{z}_i)f(t|s = 1; \mathbf{x}_i(t))] + (1 - c_i) \ln [1 - \Phi(\boldsymbol{\alpha}'\mathbf{z}_i) + \Phi(\boldsymbol{\alpha}'\mathbf{z}_i)S(t|s = 1; \mathbf{x}_i(t))],$$

For those who are observed smokers in the sample,  $c_i = 1$  and the contribution is simply the logarithm of the probability of being a smoker,  $\Phi(\boldsymbol{\alpha}'\mathbf{z}_i)$ , multiplied by the probability density function of starting at the observed starting age,  $f(\cdot)$ . For those who are observed as not starting ( $c_i = 0$ ) the contribution is the logarithm of the probability of never starting,  $1 - \Phi(\boldsymbol{\alpha}'\mathbf{z}_i)$ , plus the probability of starting after the age observed at the time of the survey,  $\Phi(\boldsymbol{\alpha}'\mathbf{z}_i)S(\cdot)$ .

Those individuals who start smoking can be linked to the rate of tax in the calendar year that they started. But this cannot be done for those respondents who had not started at the time of the HALS survey. One solution to this problem is to attribute to individuals the rate of tax in, say, the year that they were aged 16 (the modal age of starting). This is the kind of strategy adopted by Douglas and Hariharan (1994), who use the real price at age 18 and the change in real prices between the age of 15 and 18. An alternative is to treat the tax rate as a time varying covariate, as in Douglas (1998). Tax enters the likelihood function as a time varying covariate and the value of the tax rate in all of the calendar years during which the respondent is at risk of starting is used. This is the method we adopt. §

#### 4.2. Quitting

Previous models of quitting have used semi-parametric and parametric duration models to examine the effects of covariates on the number of years smoked prior to quitting. We take the sub-sample of individuals who had smoked at some point in their lives and estimate three models - the (semi-parametric) Cox proportional hazards model (Cox, 1972) and the (parametric) Weibull and generalized gamma models. ¶ All three models use time-varying covariates to model the effect of changes in the tax on cigarettes that occurred during the time an individual was a smoker on the duration of smoking.

In the Cox proportional hazards model, the hazard function at time  $t$  for individual  $i$  is defined as the product of an unspecified baseline hazard function,  $h_0(t)$ , and a proportionality factor,  $\exp(\boldsymbol{\beta}'\mathbf{x}_i(t))$ :

$$h_i(t; \mathbf{x}_i(t)) = h_0(t) \exp(\boldsymbol{\beta}'\mathbf{x}_i(t)), \quad (1)$$

where  $\mathbf{x}_i(t)$  is a vector comprising of time variant and time invariant covariates.

§Estimation of the model is carried out using STATA version 6 on a data set expanded on the age of starting for smokers and the age at the time of the survey for non-smokers. Maximum likelihood estimation is carried out using STATA's `method d0` and `method lf`.

¶This follows the logic of the split population model, which assumes that the splitting mechanism that separates those who ever smoke from those who never smoke, is independent of the hazards of starting and quitting. The nature of the quitting data in the HALS means that the duration of smoking has to be viewed as a single spell; the data does not allow analysis of multiple spells of smoking and repeated attempts to quit.

Specifying the baseline hazard function  $h_0(t)$  in Equation (1) as  $h_0(t) = pt^{p-1}$  gives the Weibull proportional hazards model, which can yield a monotonic increasing, decreasing or constant hazard rate according to the sign of  $p$ . We estimate both continuous and discrete time versions of the ‘Weibull’ model. Also we use a gamma mixture model to allow for unobservable heterogeneity. Estimation of the discrete time models exploits the fact that the data set is reshaped into longitudinal format to construct the time varying covariates. Jenkins’s (1997) estimation routine and program are used to compute discrete time ‘Weibull’ models, with and without gamma heterogeneity.

Finally, the generalized gamma model defines the hazard function as:

$$h_i(t; \mathbf{x}_i(t)) = \frac{\frac{|\kappa|}{\Gamma(\kappa^{-2})} (\kappa^{-2})^{\kappa^{-2}} \exp[\kappa^{-2}(\kappa z - e^{\kappa z})]}{\sigma t \left[ 1 - I\left(\kappa^{-2}, \kappa^{-2} \exp\left(\frac{z}{\sqrt{\kappa^{-2}}}\right)\right) \right]},$$

for  $\kappa \neq 0$ , where  $z = [\ln t - \beta' \mathbf{x}_i(t)]/\sigma$  and  $I(k, a)$  is the incomplete gamma function.||

### 4.3. Testing for misspecification

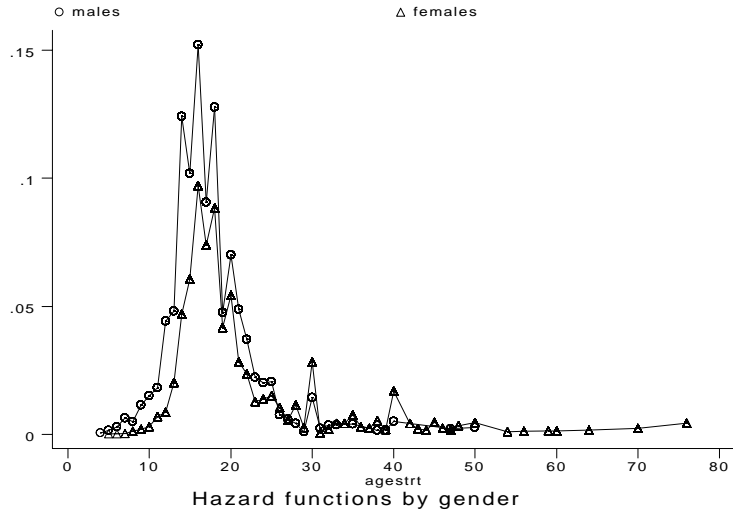
If the incorrect distribution is chosen to model either starting or quitting, parameter estimates may be biased. We use a number of diagnostic tests to test for misspecification in our models.

In the models of starting, we compare the predicted survival functions from the split population models with those obtained by (non-parametric) Kaplan-Meier estimation of the survival function for the full sample and sub-sample of smokers. We also check whether the predicted proportion of starters obtained from the split population model is close to the actual proportion observed in the data. We also use plots of the cumulative Cox-Snell residuals for the observed failures in the sample to assess the general fit of the split population model for those who fail (for a general discussion of Cox-Snell residuals see Klein and Moeschberger (1997)). A correctly fitted model should yield cumulative Cox-Snell residuals which resemble a (censored) sample from a standard exponential distribution. A plot of the non-parametric estimate of the cumulative hazard function for these data should therefore lie on a 45° line through the origin.

For the quitting data, we use the re-scaled Schoenfeld residuals (Schoenfeld, 1982) and the ‘global test’ of Grambsch and Therneau (1994) to test for non-proportionality of the hazard with respect to the covariates included in the Cox proportional hazards model. The re-scaled Schoenfeld residuals are independent of time and have an expected value of zero under the null hypothesis of proportional hazards. However, under the alternative of non-proportional hazards they will demonstrate time-dependency. We test the correlation of the residuals with a function of time - in our models this is 1 minus the Kaplan-Meier estimate of the survival function for the data. The global test for no time dependency is asymptotically distributed as a Chi-Squared variable with  $p$  degrees of freedom. In the Cox, Weibull and generalized gamma models we also use the cumulative Cox-Snell residuals and RESET-type tests as additional tests for misspecification.

We discriminate between pooled models and models split by gender using likelihood ratio tests. For the quitting data we choose an appropriate model in the light of the diagnostic tests and, to guard against over-fitting, the Akaike information criterion (AIC) (Akaike, 1974). Further we use a Wald test for  $\kappa = 1$  in the generalized gamma model (the Weibull model is a special case of the generalized gamma model when  $\kappa = 1$ ).

||The models of quitting are all estimated using STATA’s `streg` command.



**Fig. 2.** Product-limit estimates of the hazard function for starting smoking by gender.

## 5. Results

### 5.1. Starting

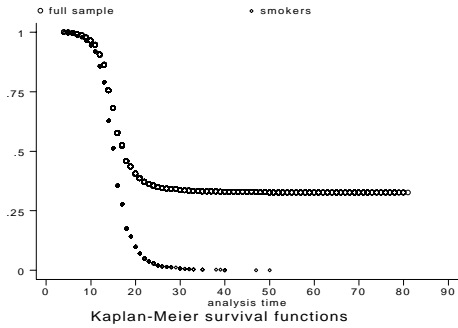
Product-limit estimates of the hazard function for starting smoking for men and women are presented in Fig. 2. The non-monotone shapes of the functions suggest that a log-logistic or log-normal specification might provide a suitable fit for the data, as was found by Douglas and Hariharan (1994) and Douglas (1998). We therefore estimate split population log-logistic/probit models and, for comparison, we also estimate log-logistic models on the full sample and the sample split by gender. LR tests of the models split by gender compared to a model estimated on the full sample indicate that the models should be analysed for men and women separately. Results for the split population models for men and women are therefore presented in Table 2.

Plots of the Kaplan-Meier estimate of the survival function for the sample which includes both starters and non-starters and the sample of starters only are presented in Figures 3 and 4 for men and women respectively. Plots of the predicted survival functions from the log-logistic/probit models which include both starters and non-starters (estimated at the mean value of the covariates) are presented in Figures 5 and 6. Finally, plots of the predicted survival functions from the split population models (plotted for the sub-sample of observed smokers only and estimated at the mean value of the covariates) are presented in Figures 7 and 8.

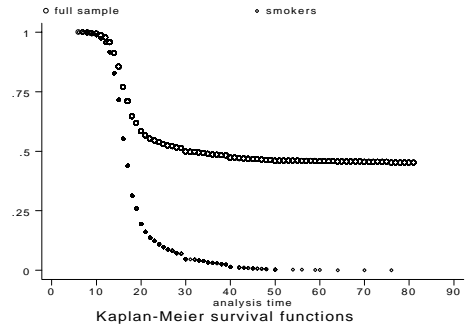
Figures 3 and 4 show that the Kaplan-Meier estimates of the survival function using the full sample yields a survival function that reaches a limit at around 0.35 for men and 0.50 for women, corresponding to the proportion of the sample who had never smoked at the

**Table 2.** Split population log-logistic/probit results for starting

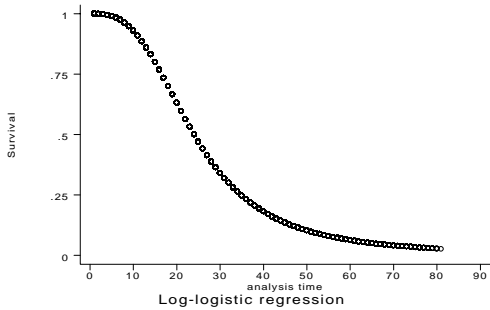
	men		women	
	duration	participation	duration	participation
NO_QUALIFICATION	-0.0453 (-3.196)	0.3289 (5.703)	-0.0402 (-2.941)	0.1601 (3.006)
A_LEVEL	0.0528 (2.121)	0.1614 (1.197)	0.0765 (2.891)	-0.1368 (-1.355)
HND	0.0300 (1.602)	-0.1623 (-1.937)	0.0196 (1.005)	-0.1914 (-2.611)
DEGREE	0.0777 (4.314)	-0.1435 (-1.752)	0.0605 (3.184)	-0.2197 (-3.106)
OTHER	-0.0110 (-0.507)	0.2683 (2.520)	0.0540 (1.785)	0.0221 (0.189)
ASIAN	0.1584 (3.001)	-0.2044 (-1.065)	0.1398 (1.652)	-0.9876 (-4.726)
AFRO-CARIBBEAN	0.1437 (2.902)	-0.0862 (-0.374)	0.0778 (1.209)	-0.4505 (-2.280)
OTHER_NON-WHITE	0.0723 (1.498)	0.1631 (0.632)	0.0246 (0.384)	-0.0169 (-0.071)
MOTHER_SMOKED	-0.0465 (-2.007)	0.4763 (4.207)	-0.0430 (-1.877)	0.5158 (5.903)
FATHER_ SMOKED	-0.0513 (-3.115)	0.4923 (6.993)	-0.0274 (-1.645)	0.2745 (4.867)
PARENTS_SMOKED	-0.0821 (-4.680)	0.4809 (6.446)	-0.0575 (-3.374)	0.5486 (9.192)
YEAR	-0.0148 (-3.964)	-	-0.0522 (-5.546)	-
YEAR <sup>2</sup> /100	0.0792 (2.703)	-	0.2626 (4.190)	-
YEAR <sup>3</sup> /1000	-0.0247 (-3.125)	-	-0.0586 (-3.718)	-
YEAR <sup>4</sup> /10000	-0.0027 (4.006)	-	0.0045 (3.434)	-
ln(TAX)	0.1635 (6.539)	-	0.0806 (3.181)	-
CONS	2.9280 (106.902)	-0.0789 (-0.944)	3.3313 (65.537)	-0.2877 (-4.345)
$\gamma$	0.1285	-	0.1342	-
Log Likelihood	-8790.46		-10297.13	
$\chi^2$ (16)	293.25		442.51	
Predicted proportion of starters	0.69		0.53	
Observed proportion of starters	0.66		0.52	
Number of observations	3737		4861	
Number of failures	2460		2508	
<i>t</i> -statistics in parentheses.				



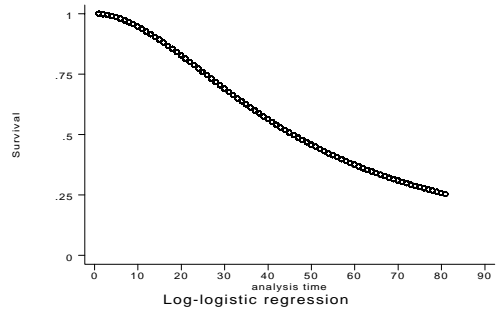
**Fig. 3.** Survival functions: Kaplan-Meier, full sample and smokers only, men



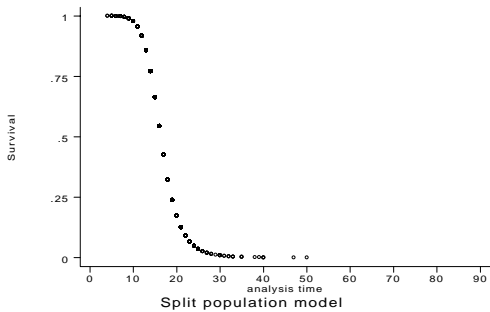
**Fig. 4.** Survival functions: Kaplan-Meier, full sample and smokers only, women



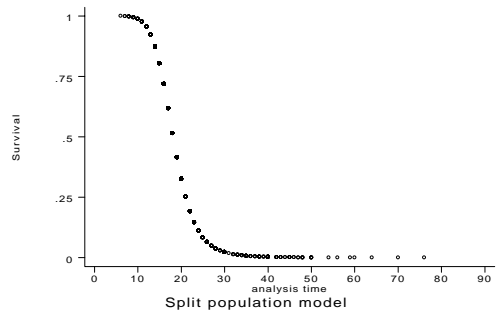
**Fig. 5.** Survival function: log-logistic, full sample, men



**Fig. 6.** Survival function: log-logistic, full sample, women



**Fig. 7.** Survival function: split population, smokers only, men



**Fig. 8.** Survival function: split population, smokers only, women

time of HALS. The Kaplan-Meier estimates of the survival functions on the sub-samples who had smoked at some point in their lives show a much steeper descent and reach zero around the late 20s for both groups, showing that virtually all smokers had started smoking by the time they were 30. The log-logistic predicted survival functions for starters and non-starters in Figures 5 and 6 should eventually reach zero, implying that all individuals in the sample will eventually 'fail' (start smoking). For men this occurs by age 80 but for women the function does not reach zero, suggesting the model is misspecified. Finally, the predicted survival functions for the split population log-logistic models in Figures 7 and 8 fit the shape of the Kaplan-Meier estimates for the subsample of smokers (shown in Figures 3 and 4) reasonably well for both men and women.

A further test of the fit of the models is given by plots of the cumulative Cox-Snell residuals. These are presented in Figures 9 to 12. Figures 9 and 10 plot the cumulative Cox-Snell residuals for the models that include both starters and non-starters, for men and women respectively. Their deviation from the 45° line indicates serious misspecification. Figures 11 and 12 plot the cumulative Cox-Snell residuals for the split population models, calculated for the sub-sample of observed starters. They suggest that the split population model for men provides a fairly good fit, as the residuals lie close to the 45° line, whereas the model is misspecified for women.

Finally, the predicted proportions of eventual starters may be calculated in the split population models and are reported at the bottom of Table 2. For men, 0.69 of the sample are predicted eventually to start, which is very close to the proportion of smokers observed in the survey (0.66). For women the figures are 0.53 and 0.52 respectively.

In Table 2 the first set of coefficients for each split population model relate to duration time and can be interpreted in terms of the qualitative effects on the age of starting. The second set of coefficients relate to the probability of never starting. For both men and women parental smoking increases the probability of starting smoking and reduces the age of starting. Measures of educational attainment suggest that those with higher levels of education are less likely to start and start later. There is a significant positive effect of an increased tax rate on the delay before starting. The estimated tax elasticity of delay is +0.16 for men and +0.08 for women.\*\*

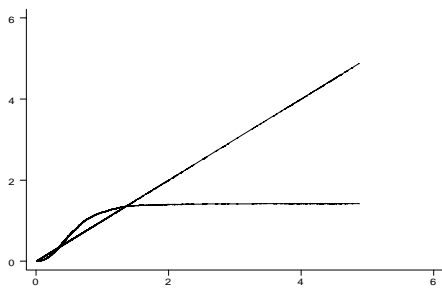
Unlike Douglas and Hariharan (1994) and Douglas (1998), who do not find a significant price effect on starting, our results do show significant effects of the tax rate on the age of starting. However, the elasticities are relatively small and do not support the evidence cited by the recent U.K. Independent Inquiry into Inequalities in Health (Department of Health, 1998a) that 'studies in the United States and Canada indicate that young people's intention to smoke and their uptake of smoking are highly price sensitive'.

### *5.1.1. Sensitivity analysis*

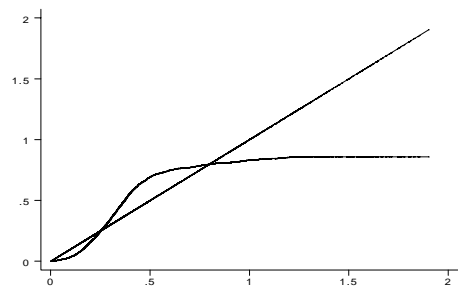
When the tax rate is treated as a time varying covariate (TVC) data is required for each year the individual is at risk, not just for the calendar year of starting. In the full sample 2.82% of individuals were born before 1905 and 17.74% before 1920. Observations prior to these calendar years have to be excluded when estimating the hazard of starting using TVCs, because of missing values for the tax data prior to 1920. We re-estimate the models

\*\*The duration models are presented in accelerated failure time format and can be interpreted as regression equations for  $\ln(\text{failure time})$ . The natural logarithm of the tax rate is used as a covariate and the coefficient can be interpreted directly as an elasticity.

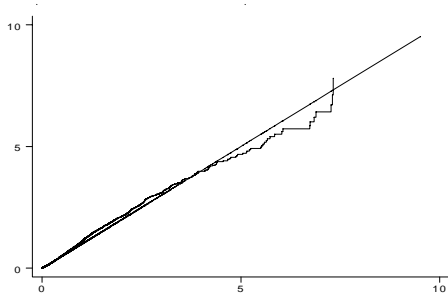




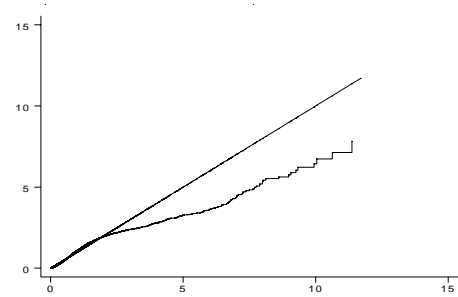
**Fig. 9.** Cumulative Cox-Snell residuals: log-logistic, full sample, men



**Fig. 10.** Cumulative Cox-Snell residuals: log-logistic, full sample, women



**Fig. 11.** Cumulative Cox-Snell residuals: split population, smokers only, men



**Fig. 12.** Cumulative Cox-Snell residuals: split population, smokers only, women

of starting using the full tax series for 1905-1985, with values for 1914-1919 set equal to the average of the 1913 and 1920 values. These new results are shown in Table 3 and show that the results are robust.

If it is the case that all non-smokers in the sample are accounted for by the splitting mechanism, and hence that none of the censored observations will eventually ‘fail’, the split population model simplifies to a ‘two-part’ specification, with a probit for starting and a log-logistic model for age of starting estimated on the sub-sample of ‘starters’. We

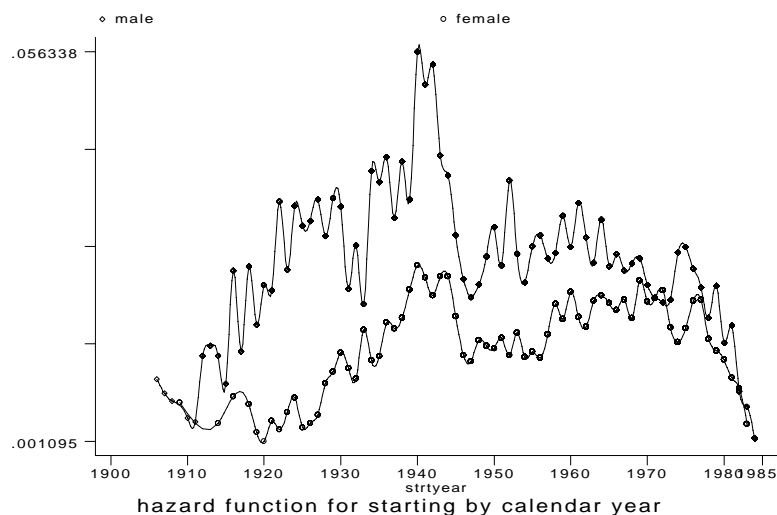
**Table 3.** Sensitivity analysis of tax elasticity in models of starting

	men	women
<i>1920-1985 data</i>		
Split population	0.164 (6.539)	0.081 (3.181)
Log-logistic	0.174 (6.638)	0.084 (2.884)
<i>1905-1985 data</i>		
Split population	0.175 (7.229)	0.107 (4.316)
Log-logistic	0.165 (6.510)	0.102 (3.551)
<i>t</i> -statistics in parentheses.		

therefore also estimate log-logistic models on the sub-samples of starters (completed spells) treating the tax rate as a TVC. Estimates of the tax elasticity in log-logistic models on the sub-samples of starters are very similar to the split population models (see Table 3).

Using retrospective data sets has been criticised as a way of analysing smoking durations by Tauras and Chaloupka (1999), who argue that asking individuals to recall events many years ago can lead to self-reporting bias in the dependent variable. To assess whether this is the case with these data, the duration data can be transformed according to the methods of Tunali and Pritchett (1997) so that the duration variable measures calendar time rather than the age of starting or quitting. The hazard functions are recalculated as a function of calendar year, with the data left truncated at the calendar year when the individual was four years old for the starting models (this being the earliest age at which an individual reported starting smoking in HALS and as such is used to define the criterion for entering the ‘risk set’ for starting smoking). Fig. 13 shows the hazard of starting smoking by calendar year using those individuals who were at risk of starting in each calendar year. There appears to be little evidence of reporting bias of the type found in the analysis of the quitting hazard using this method (see section 5.2.1). However, this does not rule out the possibility of there being other types of self-reporting bias. Fig. 13 shows a peak in the hazard, for both men and women, during the years of the Second World War and convergence of the hazard function for men and women during the 1970s and early 1980s.

With the exception of the methods used to derive the hazard function in Fig. 13, we follow Douglas and Hariharan (1994) and Douglas (1998) and use birth as the origin for the analysis of the hazard of starting. Because the estimated coefficients on  $\ln(\text{TAX})$  are elasticities, measuring the proportionate change in years before starting given a proportionate change in the tax per cigarette, they are not expected to be invariant to the choice of origin. To check the sensitivity of the results to the choice of origin, we re-estimated the log-logistic models using age 4 as the origin. This is the youngest reported age of starting in the sample, and it is reported by two individuals (in total there are 199 individuals who report starting between 4 and 10). As expected, the tax elasticities are larger in these models, but the implied effects on the number of years before starting are comparable.

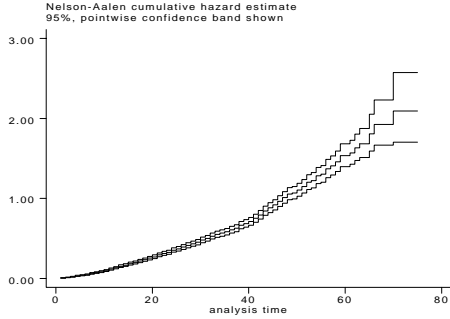


**Fig. 13.** Hazard function for starting by calendar year

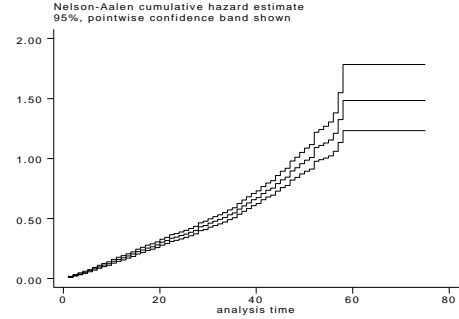
## 5.2. Quitting

The analysis of the hazard of quitting is carried out on the sub-sample of individuals who have smoked at some point in their lives.†† Those individuals who were current smokers at the time of HALS can be interpreted as ‘incomplete spells’ and defined as censored observations in the survival models. In a preliminary examination of the data we plotted Nelson-Aalen estimates of the cumulative hazard functions for men and women respectively. These are shown in Figures 14 and 15 and suggest the hazard function shows positive duration dependence. We therefore fitted three models - the Cox proportional hazards model, the Weibull model and the generalized gamma model - which allow for this shape

††Douglas (1998) applies an ordered probit specification to split the sample into three groups: those who will never start, those who start and will eventually quit and those who start and will never quit. He notes (Douglas (1998), page 55, footnote 6) that this mechanism might be modelled less restrictively as a bivariate probit. We decided to stick to a simpler splitting mechanism - based on splitting the sample by ever/never started - as we do not believe that the separation between those who will eventually quit and those who will never quit can be identified with our data. However we experimented with a bivariate probit specification, estimating sequential bivariate probit models for starting and quitting. The bivariate probit model allows for dependence between the probability of ever starting and the probability of quitting and gives us an implicit test for independence of the unobservable determinants of starting (the sample selection rule) and the hazard of quitting. The LR tests for independence ( $\chi^2(1)$ ) are 0.16 for men and 0.00 for women, providing no evidence to reject our assumption of independence. It is not possible to include the tax rate in the year of quitting in these bivariate probit models and they are estimated using all of the other regressors.



**Fig. 14.** Nelson-Aalen estimate of the cumulative hazard function for quitting, men



**Fig. 15.** Nelson-Aalen estimate of the cumulative hazard function for quitting, women

of the hazard.

Results of various tests for misspecification and splitting the model by gender are reported in Table 4 for Cox, Weibull and generalized gamma models. Table 4 serves two main purposes: it allows us to choose the most suitable model for the data on the basis of the diagnostics, and it allows us to choose whether a model split by gender is preferable to a model estimated on the full sample. The LR tests indicate that, for each specification, a model split by gender is preferred, in line with the work of Tauras and Chaloupka (1999) but in contrast to Douglas (1998).

Diagnostic tests suggest that the non-proportional hazards specification of the generalized gamma model is preferred to the proportional hazards Cox and Weibull models. The Cox models fail the re-scaled Schoenfeld residual tests on three (men) and one (women) variables, with the tax variable failing the test for proportionality for the model estimated on men and having a  $p$ -value of 0.06 in the model estimated on women. Furthermore, for men the model reports a  $p$ -value of 0.06 for Grambsch and Therneau’s global test.

Plots of the Cox-Snell residuals are shown in Figures 16 to 21 and show little difference between the models. All of the plots fit the 45° line quite well, with the generalized gamma performing slightly better for both men and women. Although the Weibull models and the generalized gamma models pass the RESET tests, both the test of  $\kappa = 1$  in the generalized gamma models and the AIC suggest that the generalized gamma model is preferred to the Weibull model and we concentrate on the generalized gamma results in our discussion. Results from the generalized gamma and Weibull models are presented in Table 5.

**Table 4.** Diagnostics for quitting

	Cox		Weibull No heterogeneity		Generalized Gamma	
	men	women	men	women	men	women
Log-likelihood	-7965.49	-6435.58	-2070.34	-2140.13	-2066.90	-2137.96
LR test of pooling samples		2892.36		81.90		71.40
Schoenfeld residuals (individual test) ( <i>p</i> )	A_LEVEL (0.03) DEGREE(0.04) ln(TAX)(0.04)	SOCIAL_CLASS_1s(0.03) ln(TAX)(0.06)	-	-	-	-
Schoenfeld residuals (global test) ( <i>p</i> )	0.06	0.14	-	-	-	-
RESET (Wald <i>p</i> )	0.88	0.24	0.910	0.28	0.79	0.513
$\kappa$ test	-	-	-	-	7.27	5.67
Akaike Information Criterion	-	-	4186.68	4326.26	4181.80	4323.92
Test of heterogeneity ( <i>p</i> values in parentheses)	-	-	21.650 (0.000)	2.332 (0.1267)	-	-
	-	-	21.152 (0.000)	2.504 (0.1135)	-	-
Notes						
LR test - model estimated on full sample against model split by gender.						
Schoenfeld residuals (individual test) - covariates which failed test for proportionality using re-scaled Schoenfeld residuals.						
Schoenfeld residuals (global test) - global test for proportionality of the hazard using method of Grambsch and Therneau (1994).						
RESET <i>p</i> -values from test of squares of predicted value of the dependent variable.						
$\kappa$ test - test for $\kappa = 1$ in the generalized gamma model. $\kappa = 1$ implies a Weibull model.						
AIC - Akaike information criterion - lowest value suggests the preferred model for males and females.						
Test of heterogeneity - LR test for Jenkins (1997) discrete time models with and without gamma heterogeneity.						

The tax elasticities of years of smoking are estimated at -0.60 for men and -0.46 for women.‡‡ As in Douglas (1998) and Tauras and Chaloupka (1999), those with higher educational qualifications have a shorter duration of smoking, as do those in higher occupational classes. Compared to the reference individual, durations of smoking are estimated to be 22% shorter for men with degrees and 29% shorter for women with degrees. Durations are estimated to be 22% longer for men with no qualifications and 20% longer for women. The occupational social class gradient is wider for women than for men. Compared to the reference individual, durations are estimated at 6% lower for men in social class 1, while for women they are estimated to be 43% lower. Similarly durations in social class 5 are estimated to be 38% higher for men and 51% higher for women.§§ Parental smoking, which is a strong predictor of starting the habit, has little effect on quitting, with only indicators that both parents smoked being significantly correlated with increased smoking duration for women. Similarly, there is little effect of the ethnic origin of the individual on the duration of smoking, other than Asian men who smoke for longer durations.

### 5.2.1. Sensitivity analysis

In order to assess the robustness of the results to the influence of recall bias, the use of continuous versus discrete time specifications of the ‘Weibull’ model, unobservable heterogeneity and the role of measures of past cigarette smoking, we carried out a sensitivity analysis of the results. Table 6 summarises the results for the estimated tax elasticity of quitting for these various models.

We experimented with different specifications of the time trend in the generalized gamma models (using polynomials of different orders to the fourth order). LR tests suggest that the fourth order time trend is jointly significant (the test statistics are 78.2 for men and 80.2 for women ( $\chi_{0.95}^2(4)$ )). Restricting the time trend to a fourth order polynomial rather than a fifth order polynomial cannot be rejected (the test statistics are 1.5 for men and 2.1 for women ( $\chi_{0.95}^2(1)$ )).

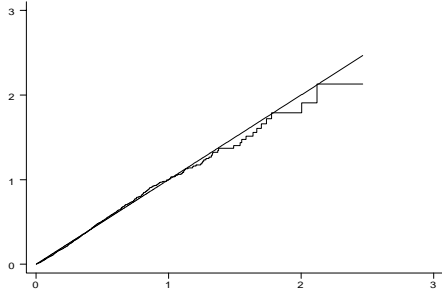
Assessment of recall bias was carried out in a similar manner to the starting models. For each individual we reconstructed the duration variable to measure calendar year of quitting, left truncated on the calendar year of starting. Results for the hazard of quitting by calendar year are presented in Fig. 22. There appears to be serious recall bias for the quitting models, due to a heaping effect, with peaks in the hazard being recorded at five and ten year intervals.¶¶ This suggests that, when questioned, respondents rounded off the ‘number of years since they quit’ (the variable EXFAGAN) to the nearest five or ten year mark.

Torelli and Trivellato (1993) propose a solution to the ‘heaping effect’ based on an explicit measurement model. They compare four methods of dealing with heaping:

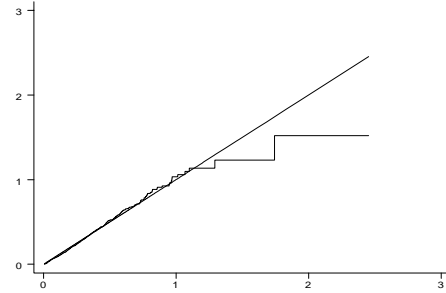
‡‡The models are again presented in accelerated failure time format and can be interpreted as regression equations for  $\ln(\text{failure time})$ . The natural logarithm of the tax rate is used as a covariate and the coefficient can be interpreted directly as an elasticity.

§§These percentages are based on  $\exp(\beta)$  which, given the semi-logarithmic specification of the model for  $\ln(T)$ , may be interpreted as elasticities.

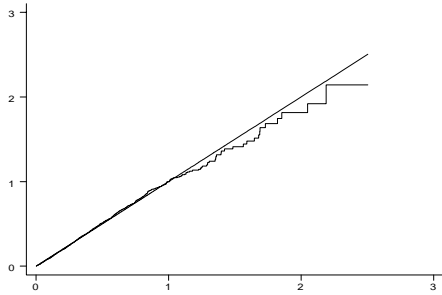
¶¶The time trend in the hazard function in Fig. 22 is informative. Prior to 1940, particularly for women, the data is too sparse to be meaningful. During the 1940s and 1950s the hazard function is relatively flat and, after the health scares of the 1960s it progressively increases.



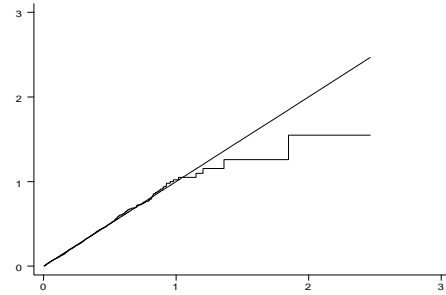
**Fig. 16.** Cumulative Cox-Snell residuals: Cox model, men



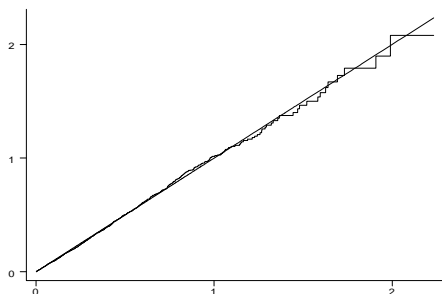
**Fig. 17.** Cumulative Cox-Snell residuals: Cox model, women



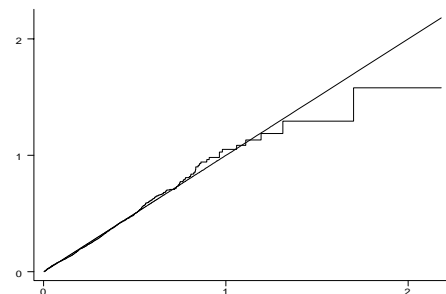
**Fig. 18.** Cumulative Cox-Snell residuals: Weibull model, men



**Fig. 19.** Cumulative Cox-Snell residuals: Weibull model, women



**Fig. 20.** Cumulative Cox-Snell residuals: generalized gamma model, men

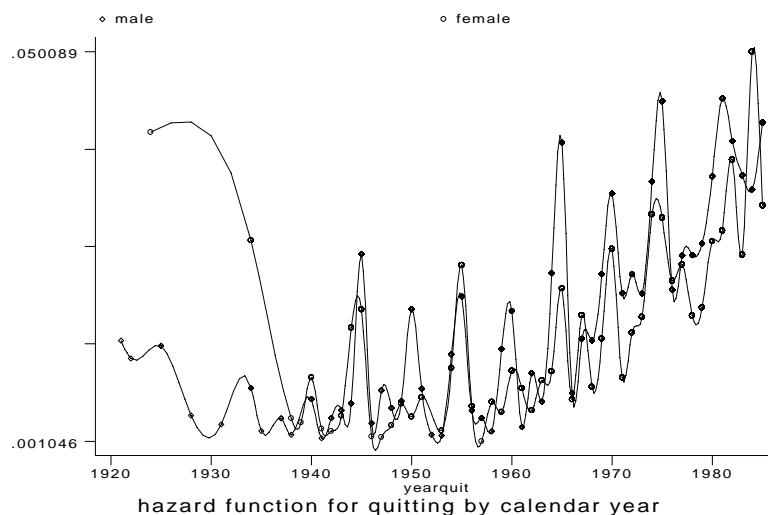


**Fig. 21.** Cumulative Cox-Snell residuals: generalized gamma model, women

**Table 5.** Accelerated failure time specifications for quitting: generalized gamma and weibull

	men		women	
	generalized gamma	Weibull	generalized gamma	Weibull
SOCIAL_CLASS_1s	-0.0664 (-0.543)	-0.0652 (-0.542)	-0.5637 (-3.588)	-0.5012 (-3.424)
SOCIAL_CLASS_2	0.0400 (0.493)	0.0411 (0.523)	-0.0332 (-0.303)	-0.0264 (-0.255)
SOCIAL_CLASS_3a	0.2082 (2.758)	0.2036 (2.791)	0.2039 (1.985)	0.2012 (2.057)
SOCIAL_CLASS_4	0.3220 (3.694)	0.3108 (3.726)	0.1607 (1.394)	0.1481 (1.342)
SOCIAL_CLASS_5n	0.3184 (2.640)	0.3120 (2.694)	0.4096 (2.469)	0.3913 (2.432)
NO_QUALIFICATION	0.1976 (2.572)	0.2049 (2.828)	0.1798 (1.888)	0.1975 (2.181)
A_LEVEL	-0.0339 (-0.235)	0.0119 (0.085)	-0.1128 (-0.655)	-0.0901 (-0.540)
HND	-0.0200 (-0.209)	-0.0178 (-0.195)	-0.1601 (-1.198)	-0.1469 (-1.176)
DEGREE	-0.2491 (-2.469)	-0.2108 (-2.273)	-0.3436 (-2.796)	-0.3163 (-2.665)
OTHER	0.3462 (3.085)	0.3531 (3.243)	-0.1905 (-1.041)	-0.1825 (-1.059)
ASIAN	0.8926 (2.442)	0.8698 (2.337)	-0.0212 (-0.033)	0.0427 (0.068)
AFRO-CARIBBEAN	-0.1771 (-0.711)	-0.2019 (-0.859)	-0.1061 (-0.225)	-0.1404 (-0.293)
OTHER_NON-WHITE	-0.0836 (-0.279)	-0.1101 (-0.410)	0.3841 (0.809)	0.3447 (0.755)
MOTHER_SMOKED	0.1738 (1.350)	0.1545 (1.257)	0.0923 (0.603)	0.0698 (0.479)
FATHER_SMOKED	-0.0542 (-0.664)	-0.0498 (-0.673)	0.2005 (1.895)	0.1883 (1.906)
PARENTS_SMOKED	0.0822 (0.901)	0.0551 (0.671)	0.2811 (2.467)	0.2400 (2.289)
YEAR	-0.0118 (-0.125)	-0.0185 (-0.182)	0.0052 (0.038)	0.0102 (0.071)
YEAR <sup>2</sup> /100	0.3281 (0.684)	0.2981 (0.589)	0.2664 (0.377)	0.2041 (0.281)
YEAR <sup>3</sup> /1000	-0.1107 (-1.137)	-0.0994 (-0.983)	-0.0948 (-0.655)	-0.0810 (-0.551)
YEAR <sup>4</sup> /10000	0.0090 (1.350)	0.0082 (1.190)	0.0072 (0.726)	0.0064 (0.637)
ln(TAX)	-0.5975 (-2.658)	-0.5285 (-2.421)	-0.4587 (-1.355)	-0.4138 (-1.235)
CONS	4.1264 (6.2027)	4.3700 (6.038)	4.1887 (4.104)	4.3277 (4.070)
ln( $\sigma$ ), ln( $p$ )	-0.2258 (-5.173)	0.3137 (11.040)	0.0280 (0.484)	0.0798 (2.586)
$\kappa$	0.7773 (9.414)		0.7792 (8.403)	
Log Likelihood	-2066.90	-2070.34	-2137.96	-2140.13
$\chi^2$ (21)	274.35	287.36	249.73	266.13
No. obs.	2480		2482	
No. failures	1176		938	
Robust <i>t</i> -statistics in parentheses.				





**Fig. 22.** Hazard function for quitting by calendar year

- (a) Re-formulating the likelihood to allow for a measurement error model, which requires specifying a parametric model of the measurement errors.
- (b) The *ad hoc* approach of adding dummy variables for the heaped observations.
- (c) Smoothing the data prior to estimation by using random draws from a uniform distribution to spread the heaped observations. This means that the results are contingent on the random numbers that are generated.
- (d) Ignoring the heaping and estimating the underlying duration model.

In their empirical application to Italian data on youth unemployment, Torelli and Trivelato find that while methods (a), (c) and (d) produce similar results, the use of dummy variables (b) produces quite different estimates. To check the sensitivity of our results to these various methods, we estimate the quitting models using all of these methods. To implement method (a), we programmed the appropriate generalized gamma model using STATA's `method lf`, assuming that the heaped observations are those where EXFAGAN is a multiple of 5 or 10. Because heaping is due to EXFAGAN the problem only relates to complete spells, that is, to those who have quit smoking. For these observations, the usual contribution to the likelihood,  $f(t_i)$ , is replaced by  $F(ut_i) - F(it_i)$ , where  $ut_i$  is the upper and  $it_i$  the lower limit of an interval of length 5 around  $t_i$ . To implement method (c), for each of the potentially heaped values (5, 10, ...), the actual observation is smoothed using pseudo random integers.† No adjustment is required for the censored observations whose durations do not depend on EXFAGAN. Method (b) simply involves adding dummy

†As the problem of heaping relates to EXFAGAN, rather than the other components of the dependent variable, we apply the smoothing method to this variable.

**Table 6.** Sensitivity analysis of the tax elasticity of quitting

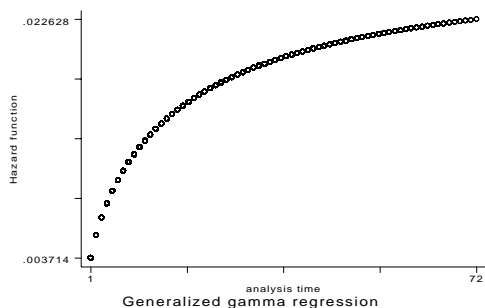
		men	women
1.	generalized gamma benchmark model	-0.60	-0.46
2.	generalized gamma heaping: dummy variables	-0.49	-0.49
3.	generalized gamma heaping: smoothed data	-0.59	-0.65
4.	generalized gamma heaping: model of measurement error	-0.49	-0.69
5.	generalized gamma with measures of actual past smoking	-0.63	-0.55
6.	generalized gamma with measures of predicted past smoking	-0.59	-0.47
7.	Weibull continuous benchmark	-0.53	-0.41
8.	'Weibull' discrete benchmark	-0.55	-0.44
9.	'Weibull' discrete with gamma heterogeneity	-0.41	-0.41

variables to indicate those observations where EXFAGAN has a potentially heaped value. Results are shown in rows 2 to 4 of table 6. Estimated tax elasticities are reduced for men and increased for women, but the sizes of the change are relatively small.

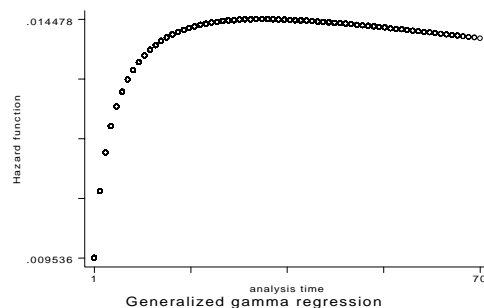
Douglas (1998) includes measures of previous cigarette consumption and of age of starting to capture addiction and lagged duration dependence in his estimates of the hazard of quitting. We have avoided these variables in our preferred models because of a concern with unobservable heterogeneity bias. However, for comparison the models are re-estimated including the age of starting smoking along with a measure of previous peak consumption. The models are estimated using both the actual values of these variables and, in an attempt to control for heterogeneity bias, the values predicted using the covariates from the models of starting. ‡ Estimated tax elasticities are reported in rows 5 and 6 of Table 6 and appear robust to these specifications; all remain negative and significant at the 5% level for men and negative but not significant for women. Estimates show little change in the value of the tax elasticity of quitting.

The specification tests presented in Table 4 suggest that the generalized gamma model is preferred to the Weibull and plots of the hazard function for the two generalized gamma models are presented in Figures 23 and 24. However, previous studies have adopted a Weibull specification and it is useful to compare the results. Row 7 of Table 6 gives the tax elasticity for the continuous time Weibull model estimated in accelerated failure time format. Very similar results are obtained for a discrete time specification - using a 'Weibull' approximation to the baseline hazard (Jenkins, 1997). These estimates are shown in row 8 of Table 6. Row 9 of Table 6 reports estimated elasticities when controlling for unobserved gamma heterogeneity using Jenkins's (1997) estimator. The estimated elasticity changes from -0.55 to -0.41 for men and from -0.44 to -0.41 for women. Table 4 shows that the unobservable heterogeneity parameter is significant in the discrete time Weibull/gamma mixture model for men but not for women.§

‡Previous peak consumption is constructed from the two variables FAGMAX and EXFAGMAX which relate to question 56(b) - 'What is the maximum number of cigarettes you have regularly smoked in a day?' and 50(c) - 'What was the maximum number of cigarettes you ever regularly smoked in a day?'.  
§We also compared estimates from proportional hazards versions of the Weibull models with those of the Cox models. There was little difference. For men, the Cox models yield a proportional shift in the baseline hazard of 2.05 (compared to 2.06 for the Weibull model), and for women the estimates are 1.53 (compared to 1.55). Elasticities are a little higher for the discrete time Weibull models compared to the continuous time Weibull models.



**Fig. 23.** Predicted hazard function for quitting: generalized gamma, men



**Fig. 24.** Predicted hazard function for quitting: generalized gamma, women

## 6. Policy implications

The goal of reducing the U.K. death rate from cancer in people under 75 by at least one fifth by the year 2010 is one of the four broad targets of the recent public health White Paper, 'Our Healthier Nation' (Department of Health, 1999). Many of these cancer deaths are attributed to smoking and this target has been linked to the Government's stated policy of reducing socio-economic inequalities in health. In July 1999, the Secretary of State for Health told the House of Commons:

'. . . tobacco smoking is the principal cause of the inequalities in health among adults in this country. Seventy per cent of deaths of working-class people, over and above what one would expect among middle class people, are the result of smoking' (Hansard, 1999)

This view is informed by the Independent Inquiry into Inequalities in Health (Department of Health, 1998a). Their review of the evidence led to the conclusion that 'smoking is an important component of differences in mortality between social classes.' Our results reinforce the evidence of the link between smoking and socioeconomic inequalities in health. The estimates in Tables 2 and 5 show that those with higher levels of education are less likely to start and, if they do, they start later. Higher educational qualifications are associated with shorter durations of smoking as are higher occupational classes.

In December 1998 the Government published a White Paper, 'Smoking Kills' (Department of Health, 1998b), to define their strategy towards smoking. The White Paper reaffirms the use of above inflation increases in tobacco taxes for health policy. Since the early 1990s successive governments have had a commitment to annual increases in the real rate of tobacco taxes, to achieve health policy objectives and encourage people to stop smoking. In the Budget of July 1997, the Chancellor of the Exchequer announced that, in future, tobacco duties would be increased on average by at least 5 per cent in real terms. This

commitment was reiterated in the 1998 and 1999 Budgets. For example, in his March 1999 statement the Chancellor of the Exchequer announced:

‘... duty on tobacco will rise by the normal escalator of 5 per cent above inflation ... a policy on cigarettes which successive British Governments have adopted for good and urgent health reasons.’ (H.M. Treasury, 1999)

Furthermore, in his November 1999 pre-Budget statement the Chancellor announced his intention that the revenue from the duty escalator would be earmarked for the National Health Service (NHS).

What are the implications of our results for this policy? Our estimated tax elasticities directly relate to the impact of above-inflation tax rises on the number of years smoked by current smokers (and indirectly on the hazards of starting and quitting). The point estimates of the tax elasticity of quitting are well defined and robust for both men and women. Point estimates are all within the range -0.40 to -0.70. This implies that the 5% real increase in tobacco duty would lead, on average, to a reduction in years of smoking of between 2% and 3.5%, where this reduction is associated with the decision to quit smoking being brought forward. Recent estimates suggest that there are around 12 million current smokers in the U.K.. Then the potential saving in total number of years smoked across the population is substantial.

Such an effect is likely to have implications for health, since smoking is recognised as a cause of mortality and as a source of differences in mortality between social classes (Department of Health (1998a), Peto et al. (1992)). Evidence also shows that quitting is beneficial for health. For example, the Whitehall I study of civil servants provides evidence on age-adjusted mortality rates over 10 years from coronary heart disease (CHD) and lung cancer (Marmot et al., 1984). For CHD, the odds ratio for deaths among ex-smokers compared to non-smokers is reported to be 1.2 while for current smokers it is 2.2. For lung cancer the odds ratio for ex-smokers is 3.0 compared to 9.4 for current smokers. The *Health and Lifestyle Survey* itself contains a follow-up of mortality rates among the original respondents. Cox et al. (1993) report the percentage mortality over the seven years from the HALS survey. Among men aged 65-74 at the time of the survey the mortality rate was 24% for those who had never smoked, 29% among ex-smokers and 39% among current smokers. For women aged 65-74 the mortality rate was 22% for non-smokers, 24% for ex-smokers and 31% for current smokers. More generally the White Paper *Smoking Kills* (Department of Health, 1998b) attributes more than 120,000 deaths per year in the U.K. to smoking, along with hundreds of deaths per year attributed to passive smoking. The costs to the NHS of smoking related illnesses are estimated at up to £1.7 billion each year (Department of Health (1998b), Parrott et al. (1998)). Taken together, these implications suggest that the impact of a continued policy of above-inflation tax rises is likely to have substantial benefits.

### **Acknowledgements**

Data from the *Health and Lifestyle Survey* were supplied by the ESRC Data Archive. Neither the original collectors of the data nor the Archive bear any responsibility for the analysis or interpretations presented here. Thanks to Brian Cox for supplying us with the estimated dates of interview for the HALS sample. Thanks to Stephen Jenkins for making available the `pgmhaz` program for STATA. Thanks to Mario Cleves for his advice on pro-

programming duration models in STATA. Thanks to two anonymous referees, Badi Baltagi, Paul Contoyannis and participants at the Ninth International Conference on Panel Data, the Bergen Workshop on Health Economics and the York Seminars in Health Econometrics for their comments. Any errors are the sole responsibility of the authors.

## References

- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transaction and automatic control AC-19*, 716–723.
- Chaloupka, F. J. (1991). Rational addictive behaviour and cigarette smoking. *Journal of Political Economy* 99, 722–742.
- Chaloupka, F. J. and M. Grossman (1996). Price, tobacco control policies and youth smoking. NBER Working paper 5740.
- Chaloupka, F. J. and H. Wechsler (1997). Price, tobacco control policies and smoking among young adults. *Journal of Health Economics* 16, 359–373.
- Cox, B. D. (1999). Personal communication.
- Cox, B. D., M. Blaxter, J. Fenner, J. Golding, and M. Gore (1987). *The Health and Lifestyle Survey* (First ed.). London: Health Promotion Research Trust.
- Cox, B. D., F. A. Huppert, and M. J. Whichelow (1993). *The Health and Lifestyle Survey: seven years on* (First ed.). Aldershot, England: Dartmouth publishing company.
- Cox, D. (1972). Regression models and life tables (with discussion). *Journal of the Royal Statistical Society* 74, 187–220.
- Department of Health (1998a). *Independent Inquiry into Inequalities in Health*. Stationery Office, London: The Stationery Office.
- Department of Health (1998b). *Smoking Kills*. Stationery Office, London: The Stationery Office.
- Department of Health (1999). *Our Healthier Nation*. Stationery Office, London: The Stationery Office.
- Dorsett, R. (1999). An econometric analysis of smoking prevalence among lone mothers. *Journal of Health Economics* 18, 429–441.
- Douglas, S. (1998). The duration of the smoking habit. *Economic Inquiry* XXXVI, 49–64.
- Douglas, S. and G. Hariharan (1994). The hazard of starting smoking: estimates from a split population duration model. *Journal of Health Economics* 13, 213–230.
- Evans, W. N. and M. C. Farrelly (1998). The compensating behaviour of smokers: taxes, tar and nicotine. *RAND Journal of Economics* 29, 578–595.
- Grambsch, P. and T. Therneau (1994). Partial hazards tests and diagnostics based on weighted residuals. *Biometrika* 81, 515–526.
- Greene, W. (1993). *Econometric analysis* (Second ed.). New York: Macmillan.
- Grossman, M. (1999). The economics of substance use and abuse. In *Paper presented at the Taipei International Conference on Health Economics, 25-27 March 1999*.
- Hansard (1999). Comment by Frank Dobson, Secretary of State for Health. House of Commons Debates 2 July 1999 c.811.

- Harris, J. E. and S. W. Chan (1999). The continuum of addiction: cigarette smoking in relation to price among Americans 15-29. *Health Economics* 8, 81–86.
- H.M.Treasury (1999). Budget speech 1999. <http://www.hm-treasury.gov.uk/budget/1999/speech.html>.
- Hsieh, C.-R. (1998). Health risk and the decision to quit smoking. *Applied Economics* 30, 795–804.
- Jenkins, S. P. (1997). sbel7 Discrete time proportional hazards regression (pgmhaz). *Stata Technical Bulletin STB-39*, 22–32. Reprinted in *Stata Technical Bulletin Reprints*, vol. 7, ed. H.J. Newton pp. 109-121, 1998. College Station TX: Stata Corporation. Code downloadable from <http://www.stata.com/support/stb/faq/>.
- Jones, A. M. (1989). A double-hurdle model of cigarette consumption. *Journal of Applied Econometrics* 4, 23–39.
- Jones, A. M. (1994). Health, addiction, social interaction and the decision to quit smoking. *Journal of Health Economics* 13, 93–110.
- Klein, J. P. and M. L. Moeschberger (1997). *Survival Analysis* (First ed.). New York: Springer-Verlag.
- Lewit, E. and D. Coate (1982). The potential for using excise taxes to reduce smoking. *Journal of Health Economics* 1, 121–145.
- Lewit, E. M., D. Coate, and M. Grossman (1981). The effect of government regulation on teenage smoking. *Journal of Law and Economics* 24, 545–569.
- Marmot, M., M. Shipley, and G. Rose (1984). Inequalities in death - specific explanation of a general pattern? *The Lancet* 1, 1003–1006.
- Parrott, S., C. Godfrey, M. Raw, R. West, and A. McNeill (1998). Guidance for commissioners on the cost-effectiveness of smoking cessation interventions. *Thorax (Supplement 5(2))* 53, S1–S38.
- Peto, R., A. Lopez, J. Boreham, M. Thun, and C. Heath (1992). Mortality from tobacco in developed countries: indirect estimation from national vital statistics. *The Lancet* 339, 1268–1278.
- Sander, W. (1995). Schooling and quitting smoking. *Review of Economics and Statistics* 129, 191–198.
- Schmidt, P. and A. D. Witte (1989). Predicting criminal recidivism using ‘split population survival time models’. *Journal of Econometrics* 40, 141–159.
- Schoenfeld, D. (1982). Partial residuals for the proportional hazards regression model. *Biometrika* 69, 239–241.
- Shmueli, A. (1996). Smoking cessation and health: a comment. *Journal of Health Economics* 15, 751–754.
- Tauras, J. A. and F. J. Chaloupka (1999). Determinants of smoking cessation: an analysis of young adult men and women. NBER Working paper 7262.

- Torelli, N. and U. Trivellato (1993). Modelling inaccuracies in job-search duration data. *Journal of Econometrics* 59, 187–211.
- Tunali, I. and J. Pritchett (1997). Cox regression with alternative concepts of waiting time: the New Orleans yellow fever epidemic of 1853. *Journal of Applied Econometrics* 12, 1–25.
- Wald, N. and A. Nicolaides-Bouman (1991). *U.K. Smoking Statistics* (Second ed.). Oxford: Oxford University Press.
- Yen, S. and A. M. Jones (1996). Individual cigarette consumption and addiction: a flexible limited dependent variable approach. *Health Economics* 5, 105–117.