

# STA442 Homework4

*SongQi Wang 1003439442*

*2019-12-03*

## Smoking

### Introduction

The age at which children first try cigarette smoking is known to be earlier for males than females, earlier in rural areas than urban areas, and to vary by ethnicity. It is likely that significant variation amongst the US states exists, and that there is variation from one school to the next.

Base on the *2014 American National Youth Tobacco Survey* ([pbrown.ca/teaching/appliedstats/data](http://pbrown.ca/teaching/appliedstats/data)), we would like to investigate the following hypotheses:

1. Geographic variation (between states) in the mean age children first try cigarettes is substantially greater than variation amongst schools. As a result, tobacco control programs should target the states with the earliest smoking ages and not concern themselves with finding particular schools where smoking is a problem.
2. First cigarette smoking has a flat hazard function, or in other words is a first order Markov process. This means two non-smoking children have the same probability of trying cigarettes within the next month, irrespective of their ages but provided the known confounders (sex, rural/urban, ethnicity) and random effects (school and state) are identical.

### Method

Children start smoking for the first time once, therefore we chose Weibull distribution to model the data, which is good for survival analysis data.

$$\begin{aligned}Y_{ijk} &\sim \text{Weibull}(\lambda_{ijk}, \kappa) \\ \lambda_{ijk} &= \exp(-\eta_{ijk}) \\ \eta_{ijk} &= X_{ijk}\beta + U_i + V_{ij} \\ U_i &\sim N(0, \sigma_U^2) \\ V_{ij} &\sim N(0, \sigma_V^2)\end{aligned}$$

where:

- $X_{ij}\beta$  is the subjects gender, ethnicity, whether they are from a rural or urban school
- $U_i$  is the school random effect.
- $V_{ij}$  is the state random effect.
- $\kappa$  is the Weibull shape parameter.

We set the following as the prior distributions:

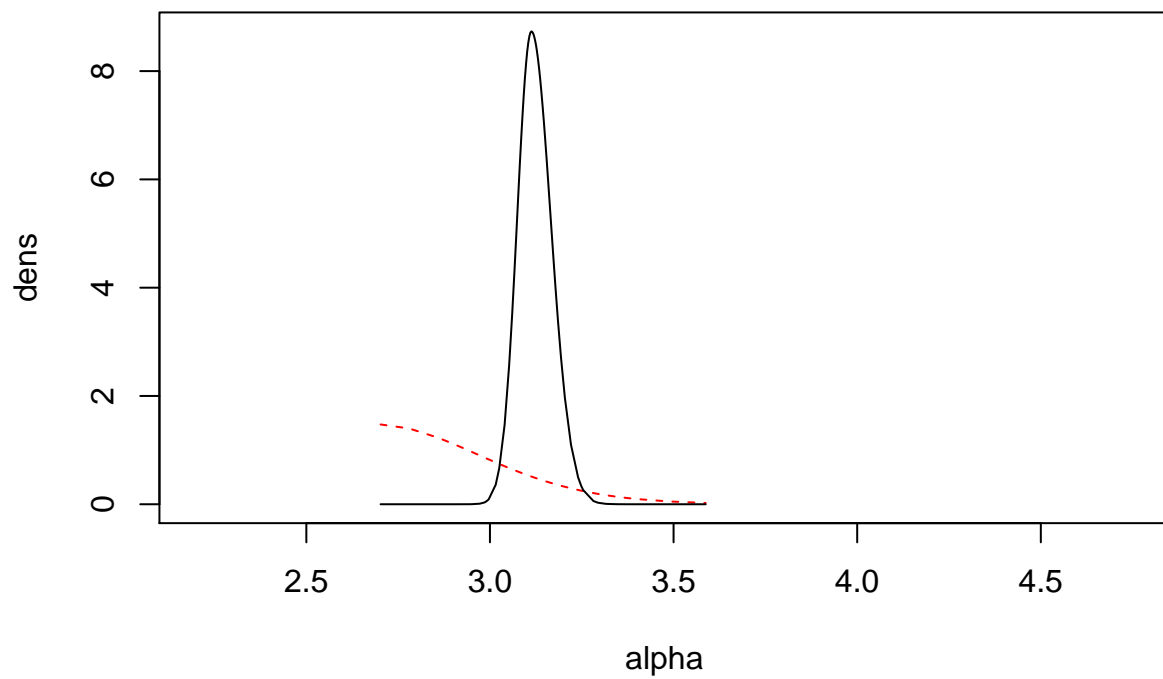
$$\begin{aligned}\kappa &\sim N(1, 0.1) \\ P(\sigma_U > 1.3) &= 1\% \\ P(\sigma_V > 0.6) &= 1\%\end{aligned}$$

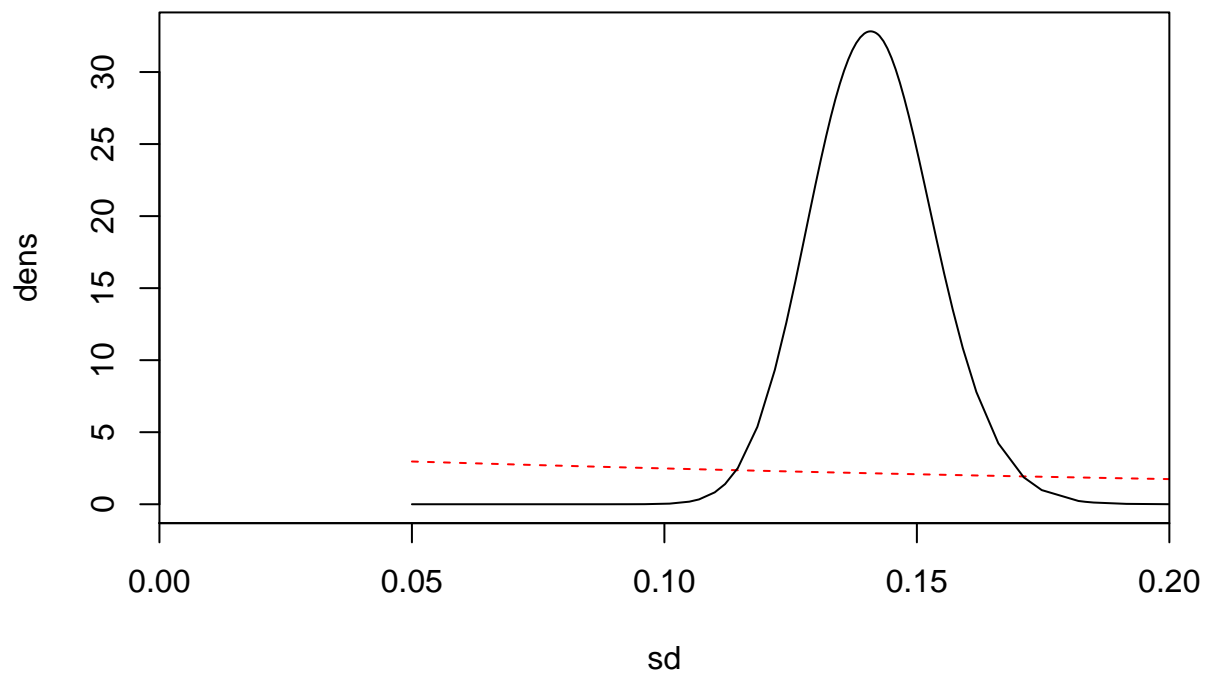
Table 1: Posterior estimates

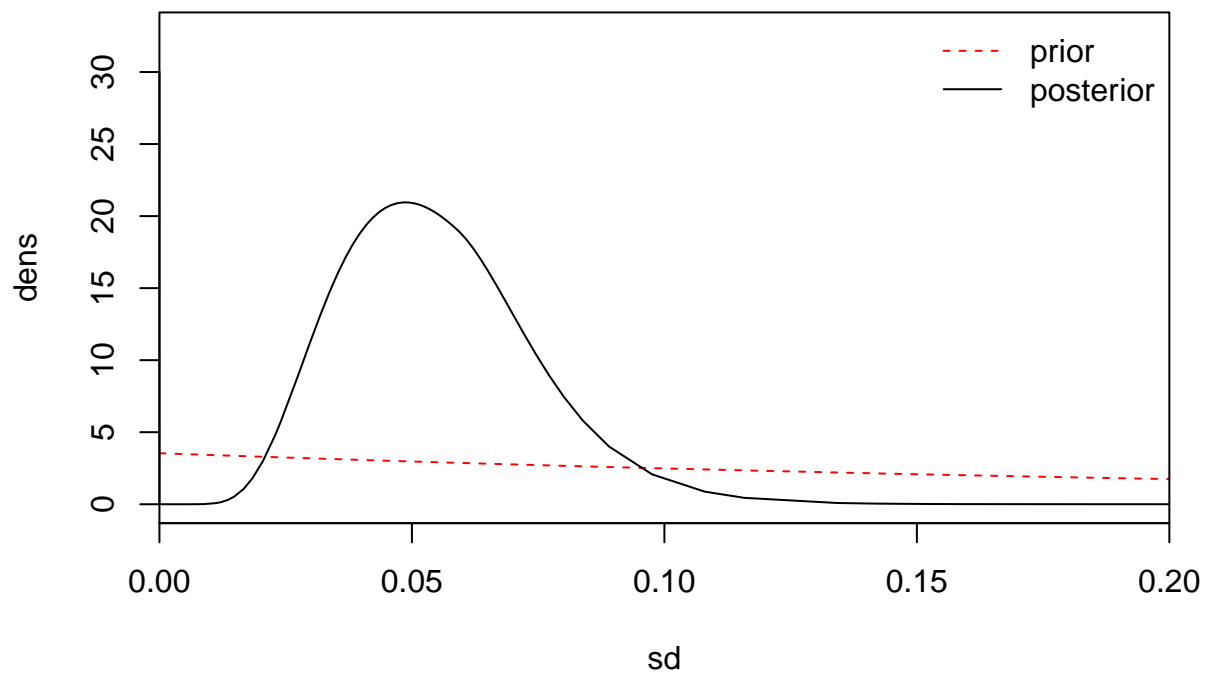
SD of School	SD of state
0.142	0.059

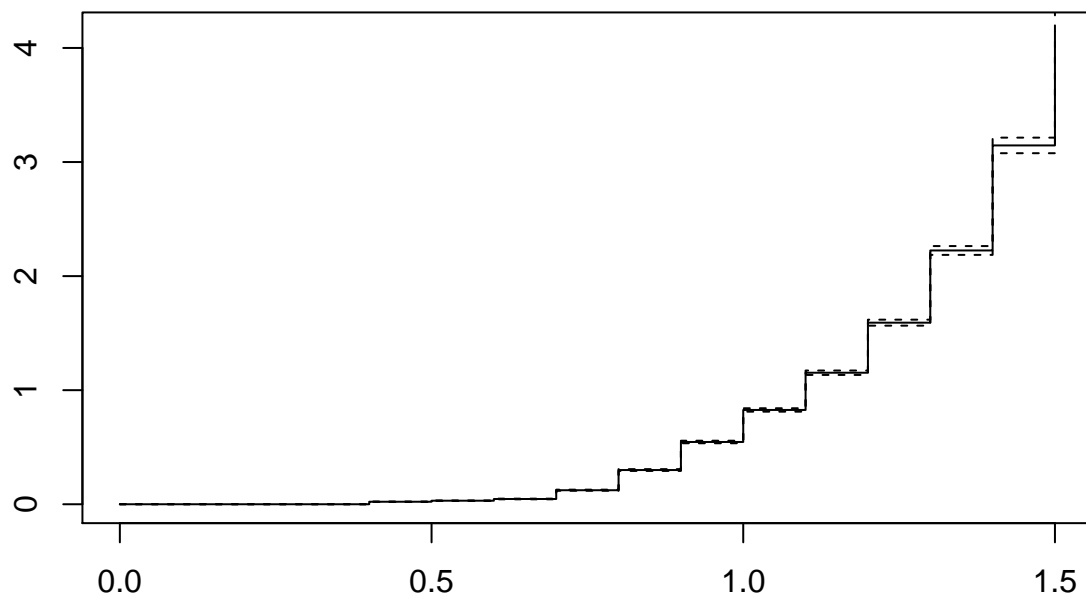
## Result

Geographic variation (between states) in the mean age children first try cigarettes is less than variation amongst schools. Tobacco control programs should actually target particular schools where smoking is a problem.









Base on the plots of prior distribution and the plot of hazard fuction, we can tell the first cigarette smoking does not have a flat hazard function. The non-smoking children with higher age have the higher probability of trying cigarettes within the next month.

# Death on the roads

## Introduction

We used the data from [www.gov.uk/government/statistical-data-sets/ras30-reportedcasualties-in-road-accidents](http://www.gov.uk/government/statistical-data-sets/ras30-reportedcasualties-in-road-accidents), the difference in accidents between the male and female, with all of the road traffic accidents in the UK from 1979 to 2015. The data below consist of all pedestrians involved in motor vehicle accidents with either fatal or slight injuries (pedestrians with moderate injuries have been removed).

We investigate that whether men are involved in accidents more than women, and the proportion of accidents which are fatal is higher for men than for women. This might be due in part to women being more reluctant than men to walk outdoors late at night or in poor weather, and could also reflect men being on average more likely to engage in risky behaviour than women.

## Method

We used conditional logistic regression to model the data. We want

$$\begin{aligned} pr(Y_i = 1|X_i) &= \lambda_i \\ \log\left(\frac{\lambda_i}{1 - \lambda_i}\right) &= \beta_0 + \sum_{p=1}^P X_{ip}\beta_p \end{aligned}$$

We have

$$\begin{aligned} pr(Y_i = 1|X_i, Z_i = 1) &= \lambda_i^* \\ \log\left(\frac{\lambda_i^*}{1 - \lambda_i^*}\right) &= \beta_0^* + \sum_{p=1}^P X_{ip}\beta_p^* \end{aligned}$$

Then we finally get:

$$\begin{aligned} \beta_p^* &= \beta_0 + \log\left(\frac{pr(Z_i = 1|Y_i = 1)}{pr(Z_i = 1|Y_i = 0)}\right) \text{ if } p = 0 \\ \beta_p^* &= \beta_p \text{ if } p \neq 0 \end{aligned}$$

where:

- $X_{ip}\beta$  is the subjects gender, and their age.
- $Y_i$  is the status of casualty.
- $Z_i$  is the strata of lightness, weather and time.

Table 2: The coefficients of conditional logistic regression

	coef	exp(coef)	se(coef)	z	Pr(> z )	sex	age
age0 - 5:sexFemale	0.0284229	1.0288306	0.0549522	0.5172285	0.6049967	Female	0
age6 - 10:sexFemale	-0.1771162	0.8376825	0.0507565	-3.4895264	0.0004839	Female	6
age11 - 15:sexFemale	-0.2498614	0.7789087	0.0471857	-5.2952744	0.0000001	Female	11
age16 - 20:sexFemale	-0.2791322	0.7564399	0.0520402	-5.3637766	0.0000001	Female	16
age21 - 25:sexFemale	-0.3691252	0.6913389	0.0633358	-5.8280613	0.0000000	Female	21
age26 - 35:sexFemale	-0.4482120	0.6387693	0.0522815	-8.5730476	0.0000000	Female	26
age36 - 45:sexFemale	-0.4482308	0.6387573	0.0516433	-8.6793515	0.0000000	Female	36
age46 - 55:sexFemale	-0.3763107	0.6863891	0.0482955	-7.7918406	0.0000000	Female	46
age56 - 65:sexFemale	-0.2370677	0.7889379	0.0403324	-5.8778460	0.0000000	Female	56
age66 - 75:sexFemale	-0.1433569	0.8664448	0.0323676	-4.4290313	0.0000095	Female	66
ageOver 75:sexFemale	-0.1256106	0.8819582	0.0272702	-4.6061492	0.0000041	Female	75
age0 - 5	0.1324083	1.1415744	0.0440170	3.0081179	0.0026287	Male	0
age6 - 10	-0.3196593	0.7263965	0.0408650	-7.8223298	0.0000000	Male	6
age11 - 15	-0.3829384	0.6818549	0.0411527	-9.3053109	0.0000000	Male	11
age16 - 20	-0.4432109	0.6419718	0.0404473	-10.9577480	0.0000000	Male	16
age21 - 25	-0.2680862	0.7648419	0.0421849	-6.3550264	0.0000000	Male	21
age 26 - 35	0.0000000	1.0000000	0.0000000	NA	NA	Male	26
age36 - 45	0.4115311	1.5091267	0.0386489	10.6479477	0.0000000	Male	36
age46 - 55	0.7682289	2.1559445	0.0389790	19.7087971	0.0000000	Male	46
age56 - 65	1.2120970	3.3605244	0.0378511	32.0227837	0.0000000	Male	56
age66 - 75	1.7972504	6.0330360	0.0363472	49.4467189	0.0000000	Male	66
ageOver 75	2.3957024	10.9759044	0.0351665	68.1244757	0.0000000	Male	75

## Result

The coefficients of conditional logistic regression are summarized in the table. The reference group is the male with age from 26 to 35. It is easy to see that generally men are involved in accidents more than women. After age 35, the proportion of accidents which are fatal is higher for men than for women, but the proportion is pretty much the same from age 0 to 35.



# Appendix