

## **Fitting Polynomials and Studying the Pattern of Prognosis of Lung Cancer in the Four Regions and Estimating the Variance of Parameters**

**Manjula S Dalabanjan**

*Research Scholar, DBIT, (Affiliated to VTU) Bengaluru-560074*

**Dr. Pratibha Agrawal**

*Prof., AMC Engg College, (Affiliated to VTU) Bengaluru-560083*

### **Abstract**

Cancer is one of the major health problem persisting worldwide. The data for prognosis of cancer is taken from the National Cancer Registry Programme. We analyse the underlying pattern of distribution of incidence rates of Lung cancer in males for the four regions Bengaluru, Chennai, Delhi and Mumbai.

We fit Model A: By observing pattern of the incidence rates of Lung cancer in males, by intuition, we divided the data into 2 groups. For Group 1, second degree equation fitted well. For Group 2, Cubic Spline model fitted well. Estimation of parameters involved in both Group 1 and Group 2 were estimated by using Least Squares Method. Expressions for the variance of parameters of both the curves were derived.

**Keywords:** incidence rates, risk, fit, Chi Square distribution, Poisson distribution, residues, cancer, regions.

### **Literature Survey**

The data for prognosis of cancer was taken from NCRP website. [a] The data gives distribution of number of people affected by all types of cancers according to age, sex, caste, region, education etc. Among these Lung Cancer, Oesophagus, Head and neck, Breast cancer and cervical cancer are some of the leading sites of cancer.

James, Vapuel, Yashin( 1986) showed that the pattern with cancer rates increasing at the decreasing rate can be modeled by a power function of the form  $\mu(x) = bx^k$  where  $\mu(x)$  is the hazard rate or incidence or mortality at age x and b and k are parameters [1].

James, Vapuel, Yashin(1985) remarked that the patterns observed may be surprisingly different from the underlying patterns on the individual level. There will be substantial

heterogeneity when theory and evidence pertaining to individuals suggest a trajectory of mortality that diverges from the observed trajectory for the populations [1].

Christopher R, Heathcote, Borek D Puza and Steven P used the data from an Australian Institute of Health Welfare (AIHW) published in AIHW(2006) a study of major causes of death in Australia during twentieth century. Generally five year age intervals were used, but they obtained a selection of numbers of deaths owing to heart or circulatory disease, cancer and other causes by single years for selected charts. The data used is of 31147 women aged 74 in 1968 and 105 in 1999 and are classified by cause of death in 3 groups.

- Group 1 being those alive in 1968 who ultimately die of heart or circulatory disease.
- Group 2 being those alive in 1968 who untimely die of cancer and
- Group 3 comprising of the remainder of 31147 women dying of other causes.

They observed that in all three cases Gompertz model is reasonable up to age 100.

Here the Gompertz model provides an acceptable description of old age mortality for all 3 causes of death separately and for the population as a whole.[2]

Balgobin, Nandaram, J Sedransk and Linda Williams Pickle (2000) used the geographical units, Health Survey Areas (HSAs), as in the Atlas (Pickle et al. 1996) which included the deaths of residents in the contiguous 48 states during 1988-1992. Here they started investigating alternative models for inference about age- specific and age – adjusted mortality rates for chronic obstructive pulmonary Disease (COPD). Their primary objective was to model mortality data for an atlas and detecting the patterns of mortality rates and identification of outliers from these patterns.[3]

Tetsuji Tonda, Kenichi Satoh et al (2011) introduced a non parametric model with time varying mixed effects for cancer. They constructed parameter estimators based on a local linear approximation and applied the proposed non parametric regression model to data collected from 47 prefectures yearly from 1975 to 2002 on large bowel cancer mortality in Japanese males. [4]

Balgobin, Nandaram, J Sedransk and Linda Williams Pickle(1999) carried out research on alternative models for estimating age specific and age adjusted mortality rates for one of the disease categories, all cancer for white males, presented in the Atlas of United states Mortality, published in 1996. They used Bayesian methods applied to four different models. [5]

Kenneth G. Manton, Eric Stallard and James W. Vaupel fitted 12 alternative models to population and mortality count for male and female cohorts. The various models evaluated are named according to the distribution of age and hazard rate assumptions employed. The likelihood ratio chi-squared test statistics for the various models were calculated. All models were estimated with 20 cohort – Specific values for the parameters  $\alpha$  and  $\beta$ . These parameters consumed 40 d.f in each analysis, leaving a residual of 180 d.f. Among these 12 models, 8 models were rejected because of heterogeneity distribution or hazard rate function within the hypothesis of homogeneity is rejected. Thus the four models which were of good fit were Gamma, Inverse Gaussian distribution, Weibull hazard rates and Gompertz distribution. [6]

P K Dhillon, B BYeole, A P Kurkure , R Dikshit and FBray (2011) wrote a paper on ‘ Trends in breast, ovarian and cervical cancer incidence in Mumbai over a 30 year period

1976-2005; an age period cohort analysis'. This study aims to quantify the recent time trends in the most commonly occurring female cancers in Mumbai women with a view to better understanding the contribution of change in life style. Breast and cervical cancer were the most frequent cancers occurring in Mumbai women together with ovarian cancer accounted for more than half of all female cancers in the study period. The rates of breast cancer among women aged 30-64 have risen gradually over the 30 – year study period with the mean increase estimated at 1.1% per year, and representing 32% of the female cancer burden in 2001-2005. In contrast cervical rates among women in the same age range decreased by 1.8% per year on average but still represents 16% of the total female cancer burden in the latest 5-year period. The age-standardized incidence of ovarian cancer among 30- to 64- year old women were reasonably stable overall, with the proportion of total female cancer incidence remaining at 7% over time.[7]

In, "A study of the uterine cervix cancer in India ", by D K Jain and Padam Singh (1996). The incidence of uterine cervix cancer varies to a great extent in the resident female population in seven cities Bengaluru, Bombay, Madras, Delhi, Ahmedabad, Poona and Nagpur in India.

The relationship between incidence and age had been studied for this cancer in which incidence increases progressively from young adult life into old age. Five year age-specific incidence rates between 30 and 69 years of age is used to estimate the parameters in each city in the mathematical model

$$I_x = bx^k$$

Where  $I_x$  is the incidence at age  $x$  and  $b$  and  $k$  are parameters.

The trends in the uterine cervix cancer based on 10 age intervals 0-29, 30-34, 35-39, 40-44, 45-49, 50-54, 55-59, 60-64, 65-69 and 70 and over, crude rate, age-adjusted rate and truncated age adjusted rate with time has been attempted by fitting an exponential model.

$$I_T = de^m$$

Where  $I_T$  is estimated annual incidence rate per 1,00,000.

T is calendar year 1981 and D and m are constants. [8]

### **Models and Notations**

The underlying pattern of distribution of incidence of Lung cancer in males was studied. There were no cases of Lung cancer observed for persons aged below 20 years. Starting from 20 to 60 years of age the data was available classified according to age, in 5 years age group.

Parabola was observed to be best fit for male young adults age ranging from 20 to 60 years. In the old adults age groups 60-64, 65-69, 70-74 and 75+ there was down turn in the incidence of Lung Cancer.

Next we observed the pattern of incidence of cancer for old adults (persons above 60 years). In the older ages the rate of incidence of cancer was observed to be deteriorating. This may be due to increase in rate of mortality in older age groups.

### Model-A

Thus let us classify the data into two groups.

Let Group 1 consist of male young adults, age ranging from 20 to 60 years.

Let Group 2 consist of male old adults, aged above 60 years.

Let  $l_{1ij}$ ,  $n_{1ij}$  and  $\lambda_{1ij}$  denote respectively, the number of males suffering from cancer, population at risk and observed rate of incidence for Group 1. Thus  $\lambda_{1ij} = \frac{l_{1ij}}{n_{1ij}}$  for age class  $j$  and in the region  $i$ .

Let  $l_{2ij}$ ,  $n_{2ij}$  and  $\lambda_{2ij}$  denote respectively, the number of males suffering from cancer, population at risk and observed rate of incidence for Group 2. Thus  $\lambda_{2ij} = \frac{l_{2ij}}{n_{2ij}}$  for age class  $j$  and in region  $i$ .

When we fitted a cubic spline model for old adults age group 60-64, 65-69, 70-74 and 75+, the cubic spline model was best fit. The estimated values obtained were almost coinciding with observed values for incidence of cancer in Group 2. Thus if  $x$  denotes age, then the estimate of incidence of Lung cancer for males in the region  $i$  and age group  $j$  is given by the following Model - A

$$\lambda_{0ij} = 0 \text{ if } x_{ij} \leq 20$$

$$\lambda_{1ij} = a_{1i}x_{ij}^2 + b_{1i}x_{ij} + c_{1i} \text{ if } 20 \leq x_{ij} \leq 60, \text{ i.e., } i = 1,2,3,4; j = 1,2,3,4,5,6,7,8$$

$$\lambda_{2ij} = a_{2i}x_{ij}^3 + b_{2i}x_{ij}^2 + c_{2i}x_{ij} + d_{2i} \text{ if } x_{ij} \geq 60, \text{ i.e., } i = 1,2,3,4; j = 9,10,11,12,$$

In the age group 0 to 20 the incidence rates of Lung cancer in males for the 4 regions Bengaluru, Chennai, Delhi and Mumbai is negligible.  $a_{1i}, b_{1i}, c_{1i}$  and  $a_{2i}, b_{2i}, c_{2i}, d_{2i}$  are the parameters which can be estimated by using least squares method for region  $i$  and for Group 1 and Group 2 respectively. Instead of using the Least squares method we can also use the Least Absolute sum of squared deviations method which give more better estimates of  $a_{1i}, b_{1i}, c_{1i}$  and  $a_{2i}, b_{2i}, c_{2i}, d_{2i}$ . However the estimates obtained by using least squares method are unbiased in nature.

We can take  $\alpha_{1ij}$  and  $\beta_{2ij}$  as normally distributed residuals of the non linear regression equations. Assume that  $\alpha_{1ij}$  be normally distributed with mean 0 and variance  $\sigma_{1i}^2$  and  $\beta_{2ij}$  be normally distributed with mean 0 and variance  $\sigma_{2i}^2$ . Thus the Model - A can be modified by adding the residuals  $\alpha_{1ij}$  and  $\beta_{2ij}$ .

$$\lambda_{1ij} = a_{1i}x_{ij}^2 + b_{1i}x_{ij} + c_{1i} + \alpha_{1ij} \text{ if } 20 \leq x_{ij} \leq 60, \text{ i.e., } i = 1,2,3,4;$$

$$j = 1,2,3,4,5,6,7,8$$

$$\lambda_{2ij} = a_{2i}x_{ij}^3 + b_{2i}x_{ij}^2 + c_{2i}x_{ij} + d_{2i} + \beta_{2ij} \text{ if } x_{ij} \geq 60, \text{ i.e., } i = 1,2,3,4;$$

$$j = 9,10,11,12,$$

Here  $\alpha_{1ij}$  and  $\beta_{2ij}$  are called as residuals. Both residuals are assumed to be normally distributed. Thus we consider  $\alpha_{1ij} \approx N(0, \sigma_{1i}^2)$  and  $\beta_{2ij} \approx N(0, \sigma_{2i}^2)$ .

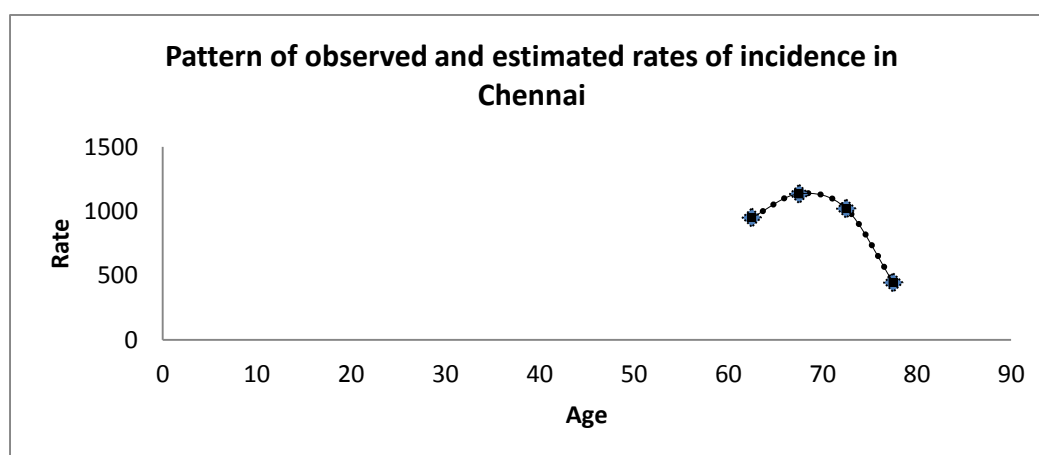
The standardized cross validation residuals for Group 1 and Group 2 are given by the following test statistic.



j	Age(mid points)	Observed Incidence Rates per 1000000	Estimated Incidence Rates	Residues	Standardised Cross Validation residual
1	27.5	8.250261871	25.60265	-17.35238813	0.09062
2	32.5	15.34074087	-1.95675	33.85291187	0.17679
3	37.5	32.08494743	16.35225	-17.35238813	0.09062
4	42.5	84.37869157	80.52965	17.29749087	0.09033
5	47.5	172.1532547	190.57545	15.73269743	0.08216
6	52.5	326.5121594	346.48965	3.849041568	0.02010
7	57.5	567.151784	548.27225	-18.42219535	0.09620

**Table 2:** For region Bengaluru and Group 2

J	Age	Observed Incidence Rates per 1000000	Estimated Incidence Rates	Residues	Standardised Cross Validation Residual
8	62.5	658.8827486	657.5684813	1.3142673	0.01737
9	67.5	828.9381394	864.4535438	-35.5154044	0.4696
10	72.5	759.850481	724.3348563	35.51562478	0.4696
11	77.5	719.6248665	720.9399188	-1.315052229	0.01737



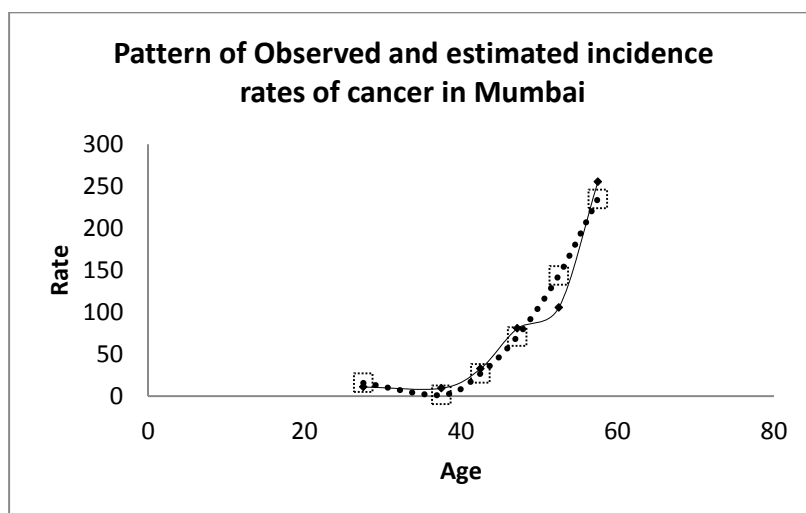
**Figure 1.3:** Group 2



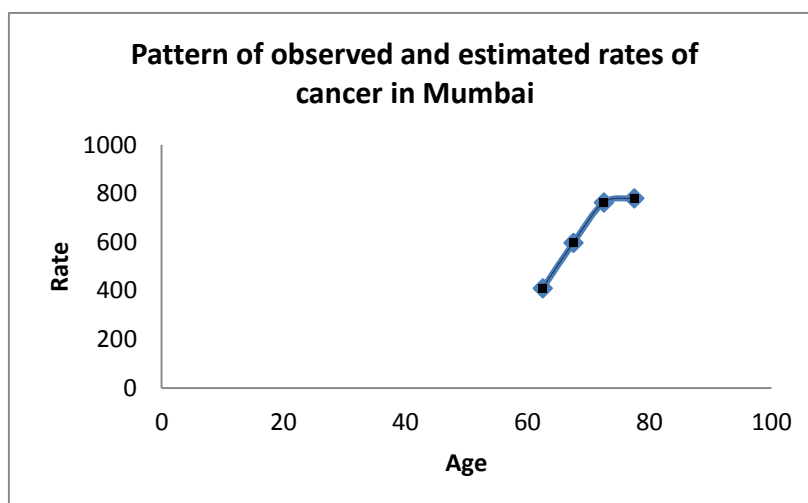


**Table 6:** For region Delhi Group 2

j	Age	Observed Incidence Rates per 1000000	Estimated Incidence Rates	Residues	Standardised Cross validation residual
8	62.5	829.2407811	829.252875	-0.012093865	0.000121
9	67.5	811.7385651	811.738625	-5.98511E-05	0.000000599
10	72.5	622.5622046	622.561375	0.000829627	0.000008314
11	77.5	620.0840452	620.071125	0.012920238	0.000129



**Figure 1.7:** Group 1



**Figure 1.8:** Group2

..... Estimated rates

\_\_\_\_\_ Observed rates

**Table 7:** For region Mumbai and Group 1

j	Age	Observed Incidence Rates per 1000000	Estimated Incidence Rates	Residue	Standardised Cross Validation Residual
1	27.5	4.316	15.861178	-11.545178	0.140426
2	32.5	6.588392899	-2.2634	4.32499	0.052606
3	37.5	9.35748396	1.462178	7.8953	0.096033
4	42.5	32.74640366	27.037678	5.708725	0.0694378
5	47.5	80.58353343	71.001478	9.582055	0.11655
6	52.5	105.4012408	143.738678	-38.3374	0.46631
7	27.5	255.1512978	234.864178	20.2871198	0.2467599

**Table 8:** For region Mumbai and Group 2

j	Age	Observed Incidence Rates per 1000000	Estimated Incidence Rates	Residue	Standardised Cross validation residual
8	62.5	409.648698	409.59125	0.05735	0.000383
9	67.5	596.8260789	596.65125	0.17482	0.0011679
10	72.5	762.448708	762.28125	0.167458	0.00111
11	77.5	780.3826662	780.48125	0.09858	0.000658

### Results and Conclusion

We find that the values of all Standardised Cross validation residual are less than 1.96 for all the four regions. Thus we can conclude that the second degree model fits well for Group 1. For Group 2 also we find that the values of all Standardised Cross validation residual are less than 1.96. Thus we can conclude that the Cubic spline model fits well for Group 2.

We can consider  $\alpha_{1ij}$  and  $\beta_{2ij}$  as normally distributed residuals of the non linear regression equations. Assume that  $\alpha_{1ij}$  be normally distributed with mean 0 and variance  $\sigma_{1i}^2$  and  $\beta_{2ij}$  be normally distributed with mean 0 and variance  $\sigma_{2i}^2$ . Thus the Model – A can be modified by adding the residuals  $\alpha_{1ij}$  and  $\beta_{2ij}$ . The residual plots for the four regions Bengaluru, Chennai, Delhi and Mumbai are shown in the following figures.

$$\lambda_{1ij} = a_{1i}x_{1ij}^2 + b_{1i}x_{1ij} + c_{1i} + \alpha_{1ij} \text{ if } 20 \leq x_{ij} \leq 60, i.e., i = 1,2,3,4;$$

$$j = 1,2,3,4,5,6,7,8$$

$$\lambda_{2ij} = a_{2i}x_{ij}^3 + b_{2i}x_{ij}^2 + c_{2i}x_{ij} + d_{2i} + \beta_{2ij} \text{ if } x_{ij} \geq 60, i.e., i = 1,2,3,4;$$

$$j = 9,10,11,12,$$

Here  $\alpha_{1ij}$  and  $\beta_{2ij}$  are called as residuals. Both residuals are assumed to be normally distributed. Thus we consider  $\alpha_{1ij} \approx N(0, \sigma_{1i}^2)$  and  $\beta_{2ij} \approx N(0, \sigma_{2i}^2)$ .

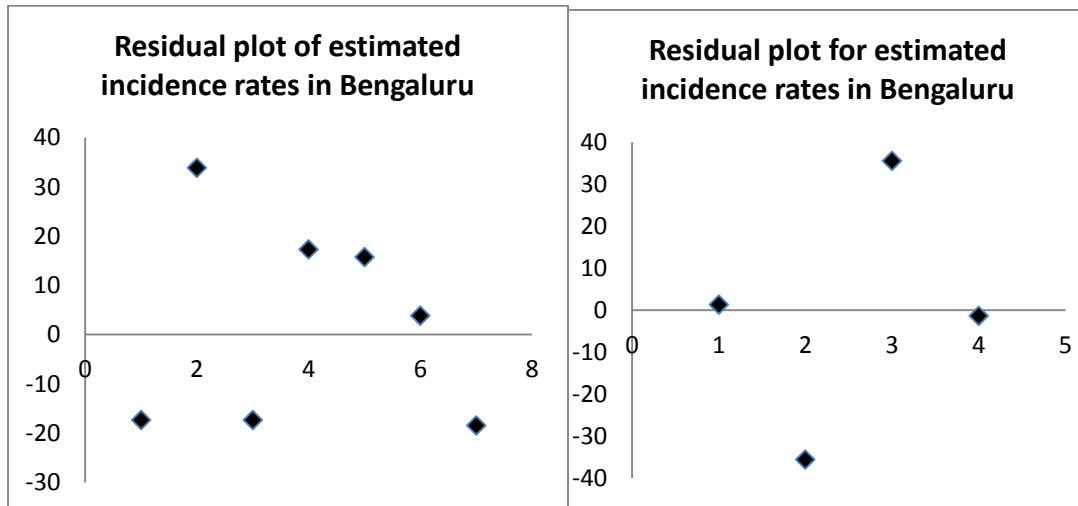


Figure 1.9: For Group 1

Figure 1.10: For Group 2

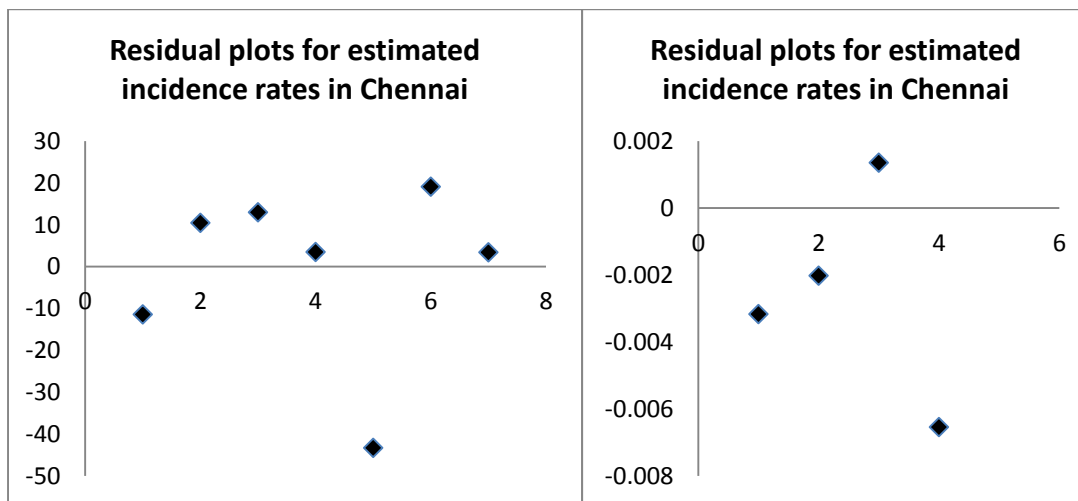


Figure 1.11: For Group 1

Figure 1.12: For Group 2

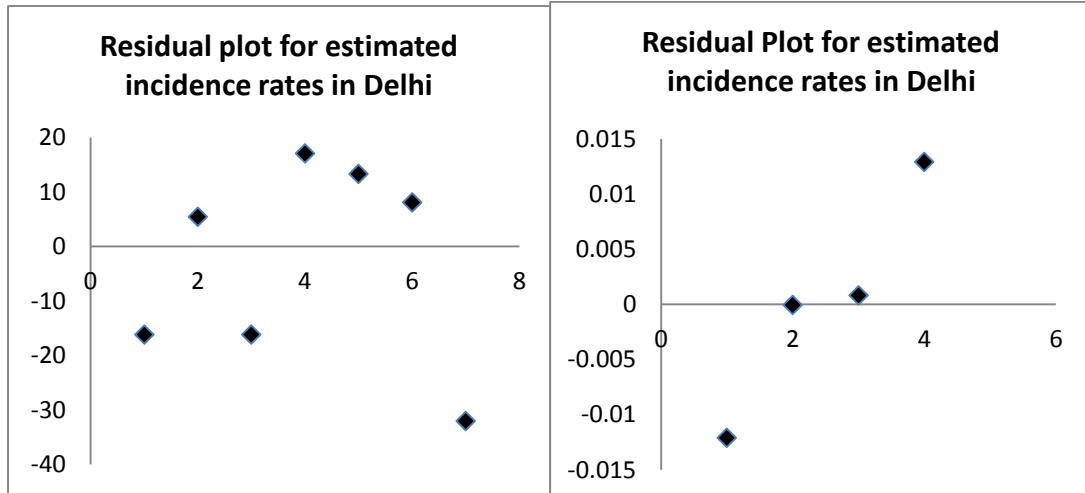


Figure 1.13: For Group 1

Figure 1.14: For Group 2

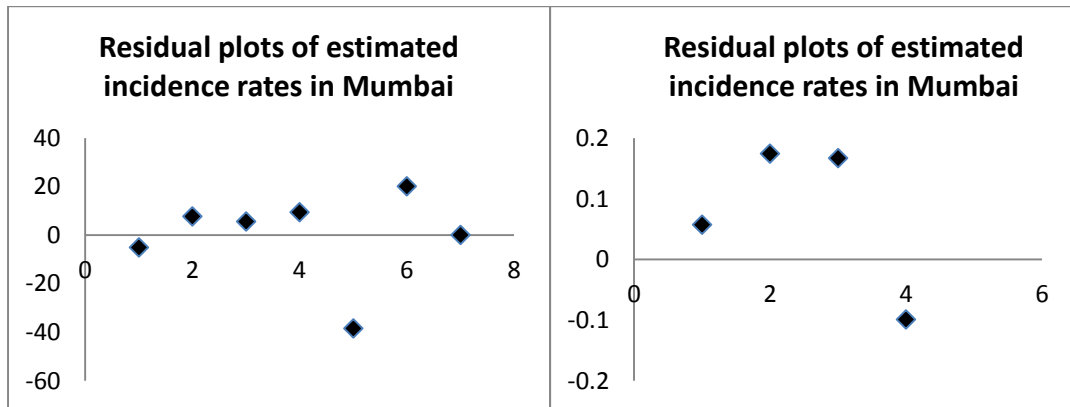


Figure 1.13: For Group 1

Figure 1.14: For Group 2

### Expressions for the variance of parameters estimated

When quadratic equation is fitted to the data the variance of parameters  $a$ ,  $b$  and  $c$  are given by the following expressions

$$y_i = ax_i^2 + bx_i + c \text{ where } \text{Var}(y_i) = \sigma^2$$

Taking  $X_i = x_i - \bar{x}$

$$\text{Var}(a) = \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

$$\text{Var}(b) = \frac{\sigma^2}{\sum_{i=1}^n X_i^2}$$

$$Var(c) = \frac{\sigma^2}{n^2} \left[ n - \left( \sum_{i=1}^n X_i^2 \right)^2 \right]$$

When the model cubic spline is fitted to the data, the variance of parameters a, b, c and d are given by the following expressions.

$$y_i = ax_i^3 + bx_i^2 + cx_i + d$$

$$Var(b) = \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \text{ Where } X_i = x_i - \bar{x}, Z_i = (x_i - \bar{x})^2 \text{ and } Var(y_i) = \sigma^2$$

$$Var(a) = \frac{\sum_{i=1}^n Z_i}{\left[ \left( \sum_{i=1}^n Z_i^2 \right)^2 - \left( \sum_{i=1}^n Z_i \right) \left( \sum_{i=1}^n Z_i^3 \right) \right]} \sigma^2$$

$$Var(c) = \frac{\sigma^2}{\sum_{i=1}^n Z_i} \left[ 1 - \frac{\left( \sum_{i=1}^n Z_i^2 \right)}{\left( \sum_{i=1}^n Z_i^2 \right)^2 - \left( \sum_{i=1}^n Z_i \right) \left( \sum_{i=1}^n Z_i^3 \right)} \right]$$

$$Var(d) = \sigma^2 \left[ \frac{1}{n} - \frac{\bar{Z}^2}{\sum_{i=1}^n (Z_i - \bar{Z})^2} \right]$$

### Determining the Odds Ratio

For region i the observed odds that a case or control will belong to group 1 is  $\frac{\lambda_{1i}}{\lambda_{2i}}$  for

control and  $\frac{\lambda'_{1i}}{\lambda'_{2i}}$  for cases.

For Bengaluru region the probability that a person having Lung cancer belongs to Group1 is around (4/10) times the probability that he belongs to Group 2. For Group 1, the odds that a person is susceptible to Lung cancer to that of he is not susceptible to Lung cancer is 0.00120. For Group2, the odds that a person is susceptible to Lung cancer to that of he is not susceptible to Lung cancer is 0.0299.

For Chennai region the probability that a person having Lung cancer belongs to Group1 is around 0.48 (1/2) times the probability that he belongs to Group 2. For Group 1, the odds that a person is susceptible to Lung cancer to that of he is not susceptible to Lung cancer is 0.00171. For Group 2, the odds that a person is susceptible to Lung cancer to that of he is not susceptible to Lung cancer is 0.00356.

For Delhi region the probability that a person having Lung cancer belongs to Group 1 is around 0.0459 (9/200) times the probability that he belongs to Group 2. For Group 1, the odds that a person is susceptible to Lung cancer to that of he is not susceptible to Lung cancer is 0.00105. For Group 2, the odds that a person is susceptible to Lung cancer to that of he is not susceptible to Lung cancer is 0.0229

For Mumbai region the probability that a person having Lung cancer belongs to Group 1 is around 0.1938 (1/5) times the probability that he belongs to Group 2. For Group 1, the odds that a person is susceptible to Lung cancer to that of he is not susceptible to Lung cancer is 0.000494. For Group 2 the odds that a person is susceptible to Lung cancer to that of he is not susceptible to Lung cancer is 0.002555.

The observed odds ratio for Group 1 in cases relative to controls is  $\hat{\Psi}$  is also called as relative risk of Lung cancer in Group 1 relative to Group 2.

$$\hat{\Psi} = \frac{\lambda_{1i} \lambda'_{2i}}{\lambda_{2i} \lambda'_{1i}}$$

Woolf's estimate of standard error of the log odds ratio is

$$S.E[\log(\Psi)] = \sqrt{\frac{1}{\lambda'_{1i}} + \frac{1}{\lambda'_{2i}} + \frac{1}{\lambda_{1i}} + \frac{1}{\lambda_{2i}}}$$

and the distribution of  $\log(\Psi)$  is approximately normal.

**Table 9:** Relative Risks In The Four Regions

Region	$\hat{\Psi}$	$S.E[\log \hat{\Psi}]$
Bengaluru	0.0405666	0.034268
Chennai	0.4805	0.02947
Delhi	0.0458837	0.03158
Mumbai	0.193436	0.04918

## References

- [1] James, Vapuel, Yashin, Heterogeneity Ruses: Some Surprising Effects of detection of population Dynamics, *The American Statistician*, Vol. 39, No. 3 (Aug 1985)
- [2] Christopher R. Heathcote, Borek D. Puza and Steven P. Roberts, *The use of aggregate data to estimate Gompertz-Type old-age mortality in heterogeneous populations*, *Australian and New Zealand Journal of Statistics*, 51(4), 481-497, (2009).

- [3] Balgobin Nandram, J Sedransk and Linda Williams Pickle, *Bayesian Analysis and Mapping of Mortality Rates for Chronic Obstructive Pulmonary Disease*, Journal of American Statistical Association vol 95, No 452, Application and Case studies (Dec 2000).
- [4] Tetsuji Tonda, Kenichi Satoh et al, *A non parametric fixed effects model for cancer mortality*, Australian and New Zealand Journal of Statistics, 53(2), 2011, 247-256.(2011)
- [5] B Nandram, J Sedransk, L Pickle, *Bayesian Analysis of Mortality Rates for U.S. Health Service Areas*, The Indian Journal of Statistics 1999, volume 61, Series B.
- [6] Kenneth G, Manton, Eric Stallard and James W Vaupel, *Alternative Models for the Heterogeneity of Mortality Risks Among the Aged*, Journal of American Statistical Association, Sept 1986, Vol. 81, No.395, Applications.
- [7] P K Dhillon , B BYeole, A P Kurkure , R Dikshit and FBray (2011), *Trends in breast, ovarian and cervical cancer incidence in Mumbai over a 30 year period 1976-2005; an age period cohort analysis*. British Journal of cancer 105,723-730.
- [8] D K Jain and Padam Singh, *A study of the uterine cervix cancer in India*, Sankhya: The Indian Journal of Statistics 1996, volume 58, series B, 118-144.
- [a] [www.ncrpindia.org](http://www.ncrpindia.org)