## 3. GENERALISED LINEAR MODELS

### 3.1 Introduction :

In classical linear model, least square theory has been used to analyse the data, where the responses are assumed to have constant variance and the systematic effects are linearly related to the mean of response variate. Then to test different hypotheses about the model parameters  $\beta$  normality is assumed. In section (2.7) it has pointed out that, classical linear models are not suitable in all situations. This is because, there are many real life situations where at least one of the assumptions of constant variance and linearity of systematic effects may not hold good. As an example consider the following real life case. Example 3.1 : An experiment was carried out on the group of 20 students, to compare three teaching methods. The students are divided into three groups and different teaching methods are applied to the various groups. At, the end of the course, a test is conducted to check the performance of these students. The marks (out of 20) obtained by them are given below.

Teaching method	1	2	3	4	5	6	7
1	19	0	5	.10i	4	17	13
11	6	12	11	9	10	13	12
111	9	8	9	; 7	8	7	1

TABLE-3.1

Below the graph of marks against teaching methods is plotted.

39

ź



The graph in Fig(3.1) shows that, the variance is not some variance stabilising transformation is constant. Hence Thus the first assumption of required. constant variance required by least square method fails. Therefore least square method can not be used to the untransformed data. In other words, classical linear model is not proper to this data. Hence it becomes necessary to search for some other models appropriate for such type of data. Alternative techniques have been developed to analyse such type of data. They are probit analysis (see e.g. Finney (1947), Prentice (1976)), logistic regression analysis (see e.g. Prentice (1976), Hosmer & Lemeshow (1989), Collett (1991)) for Binomial data; fitting log-linear models Bishop & others (1975), Good & Kurskal (1964), Gokhale & (e.g. Kullback (1978)) for the data in the form of contingency tables.

By studying similarities of many model fitting methods involving linear combinations of the parameters Nelder and Wedderburn (1972) developed the new class of models, namely, 'generalised linear models'(GLMs). This class of generalised linear models includes all the above types of models. McCullagh & Nelder (1983), Dobson (1989) have. discussed the theory associated with generalised linear model. Breslow & Day (1980)

described the applications of generalised linear models in 'cancer research'. In this chapter we discuss explicitly the theory associated with generalised linear model and the illustrations are given in the subsequent chaptersr.

Here discussion is made in the following direction.

- (1) Description of a generalised linear model with k explanatory variates,
- (2) link function,
- (3) fitting of generalised linear model,
- (4) measures of adequacy of the fitted model,
- (5) analysis of deviance (ANODEV),
- (6) model checking for generalised linear model,
- (7) method of obtaining 'robust' estimates of the model parameters in generalised linear model,
- (8) generalised linear model with varying dispersion,
- (9) fitting of generalised linear model with varying dispersion.

3.2: Generalised linear model with k explanatory variates :

The definition of generalised linear model requires the term 'one parameter natural exponential family'. But before describing this family of distributions, it is essential to explain in brief, the meaning of 'One parameter exponential family'. For this suppose,

(1)  $S \subseteq \mathbb{R}$  is the range of response variate Y, giving positive p.d.f.,

(ii)  $\theta^{-}$  is parameter of the distribution of Y,

and

۵

(iii)  $\Theta^*$  ( $\subseteq \mathbb{R}$ ) is the parameter space of  $\Theta^*$ .

<u>Definition-1</u>: <u>One parameter</u> exponential family</u> : A class of p.d.f.s(or p.m.f.s) depending on a real valued parameter  $\theta^{\#}$  of

the following form

$$f(y;\theta^*,\phi) = \exp\left\{\alpha(\phi)\left[a(y)b(\theta^*) - g(b(\theta^*)) + h(y)\right] + \xi(\phi;y)\right\} i_{\mathcal{S}}(y),$$
$$\theta^* \in \Theta^*;$$

or euivalently,

$$f(y;\theta,\phi) = \exp \left\{ \alpha(\phi) [ \alpha(y)\theta - g(\theta) + h(y) ] + \xi(\phi;y) \right\} I_{S}(y), \quad (1)$$
$$\theta^{*} \in \Theta^{*};$$

is known as one parameter exponential family, if S does not depend on  $\theta^*$ .

Here

•

- (i)  $I_S(Y) = \begin{cases} 1; & \text{if } Y \in S \\ 0; & \text{otherwise} \end{cases}$
- (ii)  $\theta = b(\theta^*)$  is a function of  $\theta^*$  only,
- (iii)  $\alpha(\phi)$  is the function of  $\phi$  alone and having positive value,
- (iv) a(Y) and h(Y) are the functions of Y only,
- (v)  $g(\theta)$  is the function of  $\theta$  only,
- (vi)  $\xi(\phi; Y)$  is the function free from  $\theta$ .

<u>Note</u>: (1) If  $\Theta^* \subseteq \mathbb{R}^k$ , distribution of Y is said to be a member of k-parameter exponential family.

(2) If the distribution of response variate Y, belongs to one parameter exponential family, the log likelihood of  $\theta$  is given by,

$$l(\theta,\phi;y) = \{\alpha(\phi)[\alpha(y)\theta - g(\theta) + h(y)] + \xi(\phi;y)\}, (2)$$

Now, we give few illustrations for understanding purpose. <u>illustration</u> 1 : <u>Normal distribution</u> :

Let the response variate Y has normal distribution  $N(\mu, \sigma^2)$ 

 $(\mu \in \mathbb{R}, \sigma \in \mathbb{R}^+)$ . Then the p.d.f. of Y is given by,

$$f(y;\mu,\sigma^{2}) = (2n\sigma^{2})^{-(4/2)} \exp\{(y-\mu)^{2}/(2\sigma^{2})\}, \quad (3)$$
  
y  $\in \mathbb{R},$ 

. . , ,

Ţ,

which gives the log likelihood function based on single observation as

$$l(\mu, \sigma^2; y) = -[\ln(2\pi\sigma^2)]/2 - (y^2 - 2\mu y + \mu^2)/(2\sigma^2). \quad (4)$$

The equation (4) can be rewritten as

$$l(\mu,\sigma^2;y) = [\mu y - (\mu^2/2) - (y^2/2)]/\sigma^2 - l_n(2\pi\sigma^2)/2.$$
 (5)

On comparing the equation (5) with the equation (2) the following findings can be obtained.

$$\phi = \sigma^{2}, \quad \alpha(\phi) = (1/\phi), \quad \theta = \mu, \quad g(\theta) = \theta^{2}/2,$$

$$a(y) = y, \quad h(y) = -y^{2}/2, \quad \xi(\phi; y) = -[\ln(2\Pi/\alpha(\phi))]/2.$$
(6)

Since the range of Y does not depend on the distribution parameter, the  $N(\mu, \sigma^2)$  distribution belongs to one parameter exponential family. Here  $\sigma^2$  is the nuisance parameter. <u>Illustration</u> 2. <u>Beta</u> <u>distribution</u> of <u>first kind</u>:

Suppose the response variate Y has beta distribution of first kind with parameters  $(\mu,\nu>0)$ . Therefore, the p.d.f. of Y is given by,  $f(y;\mu;\nu) = \{\Gamma'(\mu)\Gamma'(\nu)/\Gamma'(\mu+\nu)\} y^{(\mu-1)}(1-y)^{(\nu-1)} I_{io,ii}(y),$ 

where,

÷

t

$$I_{10,41}(Y) = \begin{cases} 0 ; if Y \in [0,1] \\ 1 ; otherwise. \end{cases}$$

Hence, the log likelihood of  $\theta$  is,

$$l(\theta, \phi; y) = ln(0) + (\mu_{\tau} 1) ln(y) + (\nu_{\tau} 1) ln(1-y), \quad (7)$$

where,

 $c = \{ \Gamma(\mu) \Gamma(\nu) / \Gamma(\mu + \nu) \}.$ 

From the equations (7) and (2), it can be seen that,

$$\phi = \alpha(\phi) = 1, \ \theta = [(\mu - 1) \ (\nu - 1)]', \ g(\theta) = 0$$

$$\underline{a}(y) = [ln(y) \ ln(1 - y)]', \ h(y) = \xi(\phi; y) = 0$$
(8)

Ξ,

1

ł

ì

As the range of Y is free from parameters of the distribution, equation (8) indicate that, beta distribution of first kind belongs to two parameter exponential family.

ł

 $R_{emark}$ : If a(y) = y, the p.d.f. of Y given in equation (1) is said to be in the canonical form, and '0' is known as canonical parameter.

Moris (1982), has named the class of distributions, having p.d.f. in the canonical form, as natural exponential family. Therefore, 'One parameter natural exponential family' can be defined as below.

<u>Definition-2</u> : <u>One parameter natural</u> exponential family :A class of p.d.f.s(or p.m.f.s) depending on a real valued parameter  $\theta^*$  of the following form

$$f(y;\theta,\phi) = \exp \left\{ \alpha(\phi) [y\theta - g(\theta) + h(y)] + \xi(\phi;y) \right\} I_{S}(y), \quad (9)$$
  
$$\theta = b(\theta^{*}) \text{ and } \theta^{*} \in \Theta^{*};$$

is known as one parameter natural exponential family, if  $S^{\frac{R}{2}}$  does not depend on  $\theta^{\frac{R}{2}}$ . The parameter  $\theta$  is called as 'natural parameter.

<u>Note</u>: If the distribution of response variate Y belongs to one parameter natural exponential family, the log likelihood of  $\theta$  is given by,

44

ŝ

$$l(\theta,\phi;y) = \{\alpha(\phi)[y\theta - g(\theta) + h(y)] + \xi(\phi;y)\}, \quad (10)$$

Note that the natural exponential family is a sub class of exponential family. Further, it can be seen that  $N(\mu, \sigma^2)$  distribution is a member of one parameter natural exponential family. On the other hand beta distribution of first kind is not a member of this class.

Below some illustrations are given to have the idea of natural exponential family more clear.

<u>illustration</u> i : <u>Gamma distribution</u> :

Let the response variate Y has gamma distribution with parameters  $(\nu, \mu)$ ,  $(\nu, \mu > 0)$  p.d.f.

$$f(y;\mu,\nu) = \{[1/\Gamma(\nu)](\nu/\mu)^{\nu}\} \exp(-\nu y/\mu) y^{\nu-1} I_{io,\infty}(y), \quad (11)$$

where

1

 $I_{to,cos}(Y) = \begin{cases} 1 ; \text{ if } Y \in [0, \infty) \\ 0 ; \text{ otherwise.} \end{cases}$ 

Thus the log likelihood function based on single observation is

 $l(\mu,\nu;y) = -\ln\Gamma(\nu) - (\nu y/\mu) + (\nu-1)\ln(y) + \nu\ln(\nu/\mu). (12)$ Rearrangement of the terms in equations (12) gives  $l(\mu,\nu;y) = \nu \left[-(y/\mu) - \ln(\mu)\right] - \ln(\Gamma(\nu)) + \nu\ln(\nu) + (\nu-1)\ln(y).$ 

(13)

Comparison of the equations (13) and (10) gives

$$\phi = (1/\nu), \ \alpha(\phi) = 1/\phi, \ \theta = -1/\mu, \quad g(\theta) = \ln(-1/\theta),$$
  
a(y) = y, h(y) = 0,  $\xi(\phi, y) = -\ln(\Gamma(\phi)) + \phi\ln(\phi) + (\phi-1)\ln(y)$ . (14)

Since the range of Y is free from the parameters  $(\mu, \nu)$ , the distribution of Y is a member of one parameter natural exponential family.

<u>Illustration</u> 2 : <u>Poisson</u> <u>distribution</u> :

Let Y be the response variate having Poisson distribution with parameter ( $\lambda$  >0) Then the p.m.f.of Y is given by

$$f(y;\lambda) = \begin{cases} exp(-\lambda)\lambda^{y}/y!; & y = 0, 1, \dots \\ 0; & \text{otherwise} \end{cases}$$
(15)

Hence the log likelihood function based on single observation y becomes

$$l(\lambda;y) = -\lambda + y \ln (\lambda) - \ln(y!)$$
(16)

Comparing the equation (16) with the equation (10) we have,

$$\phi = \alpha(\phi) = 1, \quad \theta = \ln(\lambda), \quad g(\theta) = \exp(\theta),$$
  
a(y) = y, h(y) = 0,  $\xi(\phi; y) = -\ln(y!)$  (17)

Therefore, it can be observed that, Poisson distribution belongs to one parameter natural exponential family.

<u>Illustration 3 : Binary distribution for grouped data :</u>

Let response variate Y be such that  $Y^{*}=m^{*}Y$  has Binomial distribution  $B(m^{*},p)$ . Let  $\mu$  denote the mean. Note that for grouped binary distribution  $\mu = p$ . Then the p.m.f. of Y is given by

$$f(y;m^{*},\mu) = m^{*} \left\{ m^{*}C_{m^{*}y} \mu^{m^{*}y}(1-\mu)^{(m^{*}-m^{*}y)} I_{A_{\pm}}(y) \right\}$$
(18)

Where

(i) 
$$A_{\pm} = \{0, 1/m^{\ddagger}, \dots, 1\}$$

and

È.

(ii) 
$$l_{A_{\pm}}(y) = \begin{cases} 1; & y \in A_{\pm}, \\ 0; & \text{otherwise.} \end{cases}$$

′,

\$

Thus, log-likelihood function based on single observation is  $l(m^{*}, \mu; y) = ln \left\{ \begin{array}{c} m^{*} \\ C_{m^{*}y} \end{array} \right\} + m^{*} y ln(\mu) + m^{*}(1-y) ln(1-\mu) + ln(m^{*})$ (19) The rearrangment of the terms in equation (19) gives  $l(m^{*}, \mu; y) = m^{*} [y ln [\mu/(1-\mu)] + ln(1-\mu)] + ln \left\{ \begin{array}{c} m^{*} \\ C_{m^{*}y} \end{array} \right\} + ln(m^{*})$ (20)

On comparing the equations (10) and (20), we get,  

$$\phi=1/m^*, \alpha(\phi)=1/\phi, \ \Theta = \ln(\mu/(1-\mu)), \qquad g(\Theta) = \ln(1+\exp(\Theta)),$$
  
 $a(y)=y, h(y) = 0, \ \xi(\phi, y) = \ln[m^*_{C_m^*y}] + \ln(m^*).$ 

(21)

Thus we conclude that this distribution belongs to one parameter natural exponential family.

Results of the above given distributions belonging to one parameter natural exponential family can be summarised as in the following table.

Distribution of Y	Norma I	Gamma	Poisson	Grouped binary
Range of Y	(-∞,∞)	[0,∞)	0,1,	0,1/m <sup>*</sup> ,,1
¢	a <sup>2</sup>	1/v	1	1/m <sup>*</sup>
α(φ)	1/¢	1/\$	. 1	1/¢
θ	μ	-1/µ	ln(μ)	ln{µ/(1-µ)}
g (θ)	θ <sup>2</sup> /2	ln(-1/8)	exp(0)	ln(1+exp(0))
h(y)	-y <sup>2</sup> /2	0	0	· 0 1
ζ(φ,y)	$-\ln\left[\frac{2\Pi}{\bar{\alpha}(\bar{\phi})}\right]$	-ln(Γ(φ)) +φln(φ) +(φ-1)ln(y)	-ln(y!)	$\frac{2n\left\{ {}^{m}C {}_{m} {}^{*}{}_{y} \right\}}{+ L_{n} (m^{*})}$

TABLE 3.2

A very important result which holds for one parameter natural exponential family is proved below.

<u>Result-1</u>: if the dependent variate Y comes form one parameter natural exponential family with natural parameter ' $\theta$ ' and nuisance parameter ' $\phi$ ', then

$$E(Y) = g'(\theta)$$

$$Var(Y) = g''(\theta)/\alpha(\phi)$$
(22)

<u>**Proof</u>**: Suppose Y is the single observation on the response variate with mean  $\mu$ . Assume that Y comes from one parameter natural exponential family. Hence, the log likelihood of  $\theta$  can be written as,</u>

 $\ell = \ell(\theta, \phi; y) = \{\alpha(\phi) [y\theta - g(\theta) + h(y)] + \xi(\phi; y)]\}.$  (23) It is well known that (See e.g. Kendall & Stuart (1968),

ł

$$E(\partial l/\partial \Theta) = 0, \qquad (24)$$

and

$$E(\partial l/\partial \theta)^{2} = -E(\partial^{2} l/\partial \theta^{2}), \qquad (25)$$

Differentiating equation (23) w.r.t. & twice, we get,

$$\partial l/\partial \theta = \alpha(\phi) [y - g'(\theta)]$$
 (26)

and

$$\partial^2 \boldsymbol{\ell} / \partial \boldsymbol{\Theta}^2 = \alpha(\boldsymbol{\phi}) [-g''(\boldsymbol{\Theta})], \qquad (27)$$

The equations (24) and (26) combinedly imply

 $\mu = g'(\theta); \qquad (28)$ 

and hence

$$\partial \ell / \partial \Theta = \alpha(\phi) [y - \mu]. \qquad (29)$$

From equations (25), (27) and (29), it can be seen that

 $g''(\theta) = \alpha(\phi) [var(Y)] = V.$ (30)

This result is useful in developing algorithm for fitting generalised linear model (GLM), whose definition is given below.

Definition-3 : Generalised linear model :

Let Y be a response variate with p.d.f. (or p.m.f.)  $f(.;\theta^{*})$ which belongs to the one parameter natural exponential family with  $\theta$  as a natural parameter.

Suppose  $Y_i$  (i=1,2,...,n) are n independent observations on the response variate Y. Let  $\underline{x}_1$ ,  $\underline{x}_2$ ,...,  $\underline{x}_k$  be the vectors of known values of the covariates  $X_j$  (j=1,2,...,k). Let

$$T_{i} = \beta_{o} + \sum_{j} x_{ij} \beta_{j}, \text{ for } i=1,2,\ldots,n, \qquad (31)$$

;

where

•\_

 $\beta = (\beta_0, \beta_1, \dots, \beta_k)$  is the vector unknown model parameters and

 $\underline{T} = (T_1, T_2, \dots, T_n)'$  is the linear predictor. Then a class of models of the form,

$$\underline{\mathbf{T}} = \mathbf{m}(\mathbf{E}(\underline{\mathbf{Y}})); \tag{32}$$

where m(.) is a strictly monotoic differentiable function is called the class of 'Generalised Linear Models' (GLMs).

Since '<u>T</u>' is a linear sum of the effects of explanatory variates, it is called as 'linear predictor'. Further, since the function m(.) gives the relationship between E(Y| X = x) and <u>x</u>, Nelder & Wedderburn (1972) named this function as 'link function'.

We now show that classical linear model is a particular case of generalised linear model. In classical linear model we assume normal distribution for responses. It is already shown that normal distribution is a member of one parameter natural exponential family. Further, for classical linear model, we have

$$E(\underline{Y}) = \underline{X}\beta = \underline{T}.$$

From the definition of generalised linear model, it is clear that, classical linear model is a special case of generalised linear model with 'identity' link function. By the 'identity' link function we mean the link,  $\underline{T} = \underline{\mu}$ . Thus the generalised linear model allows two extensions, namely, distribution of the response variate may be any distribution from natural exponential family, and the link function m(.) may be any strictly monotonic differtiable function.

As we have seen earlier, the functional relationship between

50

ì

Y and  $\underline{X}$  is known as link function. Along with the likelihood function, this function is also required for fitting generalised linear model to the data. Thus, link function plays a vital role in the theory of generalised linear model. Hence, it is discussed explicitly in the next section. 3.3 : Link function

Generalised linear model is fitted to the data through reparameterisation process. It should be done in such a way that the new parameters gives linear relationship between E(Y | X = x) and x. The function which relates the covariates  $(X_1, X_2, \ldots, X_k)$  and the new parameters  $\beta$  with E(Y) is known as link function. In other words, it is a function which gives relationship between linear predictor T and expected value  $\mu$  of the response variate Y. Since  $\mu$  completely depends on the behaviour of the responses, and values of the stimulus variates are fixed, the function m(.) gives link between random component and systematic components.

In classical linear models we have an identity link,  $T = \mu$ . This type of link is suitable in classical linear models, because both T and  $\mu$  can take any value on real line. Here it should be noted that such type of link is not always suitable.<sup>2</sup> Now<sup>3</sup> we illustrate this by considering two situations and suggest appropriate links in each case.

Suppose response variate Y has Poisson distribution with parameter  $\mu$ . Here identity link is not appropriate, since the linear predictor T may be negative, whereas  $\mu$  is strictly positive. So the link function m(.) must be such that m(.): $(0, \omega) \rightarrow (-\omega, \omega)$ . One such function is  $ln(\mu)$ . Thus one appropriate link in this situation is log l|ink given by  $T=ln(\mu)$ .

Secondly, assume that the response variate Y is such that

 $Y^*$  nY has binomial distribution  $B(n,\mu)$ . Here also identity link is not proper, as value of the parameter  $\mu$  is in between (0,1). Hence the link function m(.) should map from (0,1) to (- $\infty,\infty$ ). Three important functions of this type are mentioned below.

- (i) Logit function,  $T = ln \left\{-\frac{\mu}{1-\mu}\right\}$ ;
- (ii) complementary log-log function,  $T = ln(-ln(1-\mu))$ ;
- (iii) probit function,  $T = \Phi^{-1}(\mu)$ ; where  $\Phi(.)$  is the normal cumulative distribution function.

<u>Remark</u>: The link  $T = \Theta$ , is known as canonical link. Thus, 'identity link' for normal distribution, 'log link' for Poisson responses and 'logistic link' for binomial distribution are some illustrations of canonical links.

After discussing the link function explicitly, we are in a position to fit generalised linear model to the data.

## 3.4 : Fitting of GLM

Fitting the model means estimating unknown parameters in the model. Parameters in the model can be estimated by using different methods. Some of the important methods of fitting generalised linear model to the data are

(I) weighted least square method;

(II) method of maximum likelihood;

#### and

1

(III) method of obtaining robust estimates.

The third method we discuss at later stage of the chapter in section (3.9). Below we discuss the first two methods explicitly.

3.4.1 : <u>Weighted</u> <u>least</u> <u>square</u> <u>method</u> :

In least square method of estimation, estimates of the model parameters are obtained under the following two assumptions.

(1) Error components in the model are independently

distributed with mean zero and constant variance  $\sigma^2$ . (ii) The identity link holds.

in case of generalised linear models, other than classical linear model, the identity link does not hold. In other words, the systematic effects are not linearly related with the means of original responses  $Y_i$  (i=1,2,...,n). To overcome this problem, Nelder & Wedderburn (1972) defined the new variates  $Z_i$  (i=1,2,...,n) as,

$$Z_{i} = T_{i} + (Y_{i} - \mu_{i})(dT_{i}/d\mu_{i}), \qquad (1)$$

so that  $E(\underline{Z}) = X\underline{\beta}$ . This implies that, if the new dependent variable Z is considered instead of Y, the identity link is suitable for Z. Hence the generalised linear model in terms of  $\underline{Z}$ can be written as,

$$\underline{Z} = \mathbf{X} \underline{\beta} + \underline{\mathbf{e}}, \qquad (2)$$

ţ.

with

 $E(\underline{e}) = -\underline{0},$ 

and

$$Var(\underline{e}) = diag(var(\underline{e}), var(\underline{e}), \dots, var(\underline{e})), \quad (3)$$

Consider

$$\operatorname{var}(\mathbf{e}_{i}) = \operatorname{var}\left[(Y_{i}-\mu_{i})(dT_{i}/d\mu_{i})\right]$$
$$= (dT_{i}/d\mu_{i})^{2} \operatorname{var}(Y_{i})$$
$$= \left\{V_{ii}(\mu_{i})/(\alpha(\phi))\right\}(dT_{i}/d\mu_{i})^{2}$$
$$= (\alpha(\phi)W_{ii})^{-1}(\operatorname{say}),$$

where

$$W_{ii} = \left\{ \left[ V_{ii} (\mu_i) \right] (dT_i / d\mu_i)^2 \right\}^{-1}.$$
 (4)

;

лл Э

Define

$$W = diag(W_{11}, W_{22}, ..., W_{nn}).$$
 (5)

Premultiplying by  $W^{1/2}$  to the model (2), we have

$$\underline{Z}^{*} = X^{*} (\underline{3} + \underline{9}^{*}), \qquad (6)$$

$$\underline{Z}^{*} = W^{(1/2)} \underline{Z}, \qquad (7)$$

$$X^{*} = W^{(1/2)} X, \qquad (7)$$

and

with

Since  $\alpha(\phi)$  is assumed to be constant, for the model (6), one can obtain least square estimates of the model parameters  $\beta$  by using well known normal equations,

or equivaletly,

$$(\mathbf{X'W} \mathbf{X})_{\beta} = \mathbf{X'W} \underline{Z}.$$
 (7)

Since in generalised linear model, the newly defined variate Z and the weight matrix W both depend on the estimates of model parameters  $\partial_i$ , it is necessary to use the weighted least square method iteratively. Here it is essential to note that, as the newly defined dependent variate Z is unobservable, it is not possible to obtain weighted least square estimates of  $\partial_i$  without making any assumption about the distribution of response variate. Hence, for fitting generalised linear model, this method can be used only after making a valid distributional assumption for the response variate Y. Once we assume a particular distribution

ز

from one parameter natural exponential family, this method can be used.

<u>Note</u> :- In case of classical linear model we have  $E(Y) = \mu$  and  $V(\mu) = 1.$  (8) Therefore we have from equations (1) and (4)

$$Z = Y \text{ and } W = 1 \tag{9}$$

Hence for n observations equations' (7) reduces to

$$(\mathbf{X'X})\boldsymbol{\beta} = \mathbf{X'Y},$$
 (10)

which are the normal equations given in equation (2.3-6). Thus in case of classical linear model, weighted least square method is equivalent to the usual least square method.

Now we discuss the second method.

3.4.2 : <u>Method of maximum likelihood</u> :

ì

In the class of generalised linear models, distribution of the response variate is a member of one parameter natural exponential family. Hence, the method of maximum likelihood can be used to estimate the model parameters (3.

Suppose  $Y_i$  (i=1,2,...,n) are *m* independent responses with p.d.f. (or p.m.f.) of  $Y_i$  is given by,

$$f(y_i;\theta_i,\phi) = \exp\{\alpha(\phi)[y_i\theta_i - g(\theta_i) + h(y_i)] + \xi(\phi;y_i)\}\}.$$

Therefore, the log likelihood of  $\underline{\Theta}$  based on n observations is given by,

$$\mathcal{E}(\underline{\Theta}, \phi; \underline{Y}) = \sum_{i} \{ \alpha(\phi) [ y_i \Theta_i - g(\Theta_i) + h(y_i) ] + \beta(\phi_i, y_i) ] \}. (11)$$

In order to obtain maximum likelihood estimates of  $\beta$ , we differentiate the expression (11) and set the results equal to zero. Differentiation of expression (11) can be obtained by using chain rule. According to this rule, we can write

$$\frac{\partial \ell}{\partial \beta_j} = \sum_i \left( \frac{\partial \ell}{\partial \Theta_i} \right) \left( \frac{\partial \Theta_i}{\partial \Theta_i} \right) \left( \frac{\partial \mu_i}{\partial \Phi_i} \right) \left( \frac{\partial \mu_i}{\partial \Phi_i} \right) \left( \frac{\partial T_i}{\partial \Phi_i} \right), \quad (12)$$

From equation (11) it can be observed that,

$$(\partial l/\partial \theta_{i}) = \alpha(\phi) [y_{i} - g'(\theta_{i})]$$
(13)

Further, from the result (3.2-22) we have,

$$\left. \begin{array}{c} g'(\Theta_{i}) = \mu_{i} \\ g''(\Theta_{i}) = V_{ii}(\mu_{i}) \end{array} \right\}$$
 (14)

•

.

Also, we have

$$F_{i} = \sum_{j} x_{ij} \beta_{j},$$
  
 $i = 1, 2, ..., n;$   
 $x_{i0} = 1;$  for all i.

Hence,

$$\partial T_i / \partial \beta_j = x_j ; j=0,1,\ldots,k.$$
 (15)

Using equations (13) to (15), in (12) we get the estimating equations as

$$T_{j}^{*} = \sum_{i} \left\{ \begin{array}{c} \alpha(\phi) \left[ y_{i} - \mu_{i} \right] & d\mu_{i} \\ ------- & ----- & ----- \\ V_{i}(\mu_{i}) & dT_{i} \end{array} \right\} = 0$$
(16)  
$$j = 0, 1..., k.$$

Alternatively, the above equations (16) can be written in the matrix form as,

$$D'C^{-1}(\underline{y}-\underline{\mu}) = \underline{0}$$
, (17)

where

(i)  $D = (D_{ij})$ , with  $D_{ij} = (d\mu_i/d\beta_j)$  for every i and j value, (ii) C is the variance covariance matrix of <u>Y</u>.

;

After solving equations (17) we obtain estimates' of  $\beta_{j}(j=0,1,\ldots,k)$ .

Since Y's (i=1,2,...,n) are independent,

$$C = diag(var(Y_1), var(Y_2), \dots, var(Y_n)).$$
(18)

It can be observed that in case of classical linear model, the estimating equations (17) reduces to the normal equations given in (2.3-6). This can be justified as follows.

Justificatin : For classical linear model, we have,

 $\frac{\mathbf{T}}{\mathbf{T}} = \boldsymbol{\mu}, \qquad (19)$  $\mathbf{C} = \sigma^2 \mathbf{I}_{\mathbf{n}},$ 

and,

where  $\underline{T} = X\beta$  and  $\mu = E(\underline{Y})$ . Equation (19) implies,  $\mathbf{D} = \mathbf{X}$ , and  $\mu = \mathbf{X}\beta$ . Hence for classical linear model, the estimating equations (17) become,

 $\mathbf{X}^{*}(\underline{\mathbf{y}}-\mathbf{X}\beta) = \underline{\mathbf{0}}.$ 

i.e.

$$(\mathbf{X'X})_{\boldsymbol{\beta}} = \mathbf{X'Y}, \qquad (20)$$

which are the normal equations as given in equations (2.3-6). Since the normal equations in (20) are linear equations in  $\beta$ , the estimates of  $\beta$  can be obtained by solving them.

For non normal distributions, the equations (17) are non linear equations in  $\beta$ . Therefore, it may not be possible to obtain the maximum likelihood estimates explicitly. However by using appropriate numerical technique, the estimates of  $\beta$  can be obtained.

For the class of generalised linear models defined by Nelder & Wedderburn (1972), the approach of Newton-Raphson method is

ļ

used to compute maximum likelihood estimates. Neider & Wedderburn (1972) introduced this method. Later McCullagh & Neider (1983) reproduced the same with more explanation. We also use the Newton-Raphson method as a numerical technique to solve the equations (17). According to this method  $m^{th}$  approximation of the estimates of  $\beta$  is given by

$$\hat{\boldsymbol{\beta}}^{(m)} = \hat{\boldsymbol{\beta}}^{(m-1)} - \left[\underline{\boldsymbol{T}}^{*(m-1)}\right]_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}^{(m-1)}} \cdot \left[ -\frac{\partial^2}{\partial \beta_j} \frac{\partial}{\partial \beta_l} \right]_{\hat{\boldsymbol{\beta}}=\hat{\boldsymbol{\beta}}^{(m-1)}}^{-1} ; \quad (21)$$

where  $\left[-\frac{\partial^2}{\partial \beta_j} \frac{l}{\partial \beta_1}\right]$  is the matrix of second order derivatives of l. Hence to solve the equations (21) it is essential to obtain the expressions for  $\left[-\frac{\partial^2}{\partial \beta_j} \frac{l}{\partial \beta_1}\right]$ . For this consider,

$$-\frac{\partial^{2}}{\partial \beta_{j}}\frac{\partial}{\partial \beta_{l}} = -\frac{\partial}{\partial \beta_{l}}\left[\frac{\partial \ell}{\partial \beta_{l}}\left[\frac{\partial \ell}{\partial T_{i}}, \frac{\partial T_{i}}{\partial \beta_{j}}\right]\right]$$

$$= \sum_{i} -\frac{\partial}{\partial \beta_{l}}\left[\frac{\partial \ell}{\partial T_{i}}, \frac{\partial T_{i}}{\partial \beta_{j}}\right]$$

$$= \sum_{i} -\frac{\partial}{\partial \beta_{l}}\left[\frac{\partial \ell}{\partial T_{i}}, x_{ij}\right]$$

$$= \sum_{i} -\frac{\partial}{\partial T_{i}}\left[\frac{\partial \ell}{\partial T_{i}}, \frac{\partial T_{i}}{\partial \beta_{l}}\right] x_{ij}$$

$$= \sum_{i} -\frac{\partial^{2}}{\partial T_{i}^{2}}, (x_{ij}x_{il}), \qquad (22)$$

Consider,

$$\frac{\partial^{2} \ell}{\partial T_{i}^{2}} = -\frac{\partial}{\partial \bar{T}_{i}} \begin{bmatrix} \partial \ell & d\theta_{i} \\ \bar{\partial}\bar{\theta}_{i} & -\bar{d}\bar{T}_{i} \end{bmatrix}$$

$$= \frac{\partial \ell}{\partial \bar{\theta}_{i}} - \frac{d^{2} \theta_{i}}{dT_{i}^{2}} + -\frac{d\theta_{i}}{d\bar{T}_{i}} - \frac{\partial}{\partial \bar{\theta}_{i}} \begin{bmatrix} \partial \ell \\ \bar{\partial}\bar{T}_{i} \end{bmatrix}$$

58

ř

$$= \frac{\partial \ell}{\partial \theta_{i}} \cdot \frac{d^{2} \theta_{i}}{dT_{i}^{2}} + \frac{d \theta_{i}}{d\overline{T}_{i}} \cdot \frac{\partial}{\partial \overline{\theta}_{i}} \left[ \frac{\partial \ell}{\partial \overline{\theta}_{i}} \cdot \frac{\partial \theta_{i}}{\partial \overline{T}_{i}} \right]$$
$$= \frac{\partial \ell}{\partial \overline{\theta}_{i}} \cdot \frac{d^{2} \theta_{i}}{d\overline{T}_{i}^{2}} + \left[ \frac{d \theta_{i}}{d\overline{T}_{i}} \right]^{2} \cdot \left[ \frac{\partial^{2} \ell}{\partial \overline{\theta}_{i}^{2}} \right].$$
(23)

Since the response variate is assumed to have a distribution from one parameter natural exponential family, we get from the equations (3.2-27) and (3.2-28),

$$\partial^2 \ell / \partial \Theta_i^2 = -\alpha(\phi) V_{ii}(\mu_i). \qquad (24)$$

$$d\mu_i / d\theta_i = V_{ii} (\mu_i), \qquad (25)$$

, ' ' †

Using equations (3.2-29) and (24), equation (23) can be rewritten as,

$$-\frac{\partial^{2} t}{\partial T_{i}^{2}} = \alpha(\phi) \left\{ \begin{pmatrix} y_{i}^{-} \mu_{i} \end{pmatrix} - \frac{d^{2} \theta_{i}}{\partial T_{i}^{2}} - V_{ii}(\mu_{i}) \begin{bmatrix} \frac{d \theta_{i}}{\partial \mu_{i}} & \frac{d \mu_{i}}{\partial T_{i}} \end{bmatrix}^{2} \right\}.$$
 (26)

Sometimes, Fisher's scoring method is simpler than Newton-Raphson method. In Fisher's method of scoring, matrix of second order derivatives in equation (21) is replaced by its expectation. Now taking expected value equation (26) can be written as

$$E\left[-\frac{\partial^{2}}{\partial \beta_{j}}\frac{\partial}{\partial \beta_{l}}\right] = -\sum_{i}\left\{-\frac{\alpha(\phi)}{V_{ii}}\frac{\partial}{\partial \mu_{i}}\right\}\cdot\left[-\frac{d\mu_{i}}{d\overline{T}_{i}}\right]^{2}x_{ij}x_{il}$$
$$= -\sum_{i}\alpha(\phi)W_{ii}(x_{ij}x_{il}). \qquad (27)$$

Hence by proper changes in equation (21), and using equation (27) we get,

$$\hat{\beta}^{(m)} = \hat{\beta}^{(m-1)} + \left[\underline{T}^{*(m-1)}\right]_{\hat{\beta}=\hat{\beta}^{(m-1)}} \cdot \left[\alpha(\phi)\sum_{i} \{W_{ii} x_{ij} x_{ii}\}\right]_{\hat{\beta}=\hat{\beta}^{(m-1)}}^{-1}$$

$$= \hat{\beta}^{(m-1)} + \left[\sum_{i} \left\{\frac{\alpha(\phi)(x_{ij}(y_{i}-\mu_{ij}))}{V_{ii}}\right\}_{V_{ii}}^{-1}\right] \cdot \left[d\mu_{i}/dT_{i}\right]_{\hat{\beta}=\hat{\beta}^{(m-1)}}^{-1}$$

$$\cdot \left[\alpha(\phi)\sum_{i} \{W_{ii} x_{ij} x_{ii}\}\right]_{\hat{\beta}=\hat{\beta}^{(m-1)}}^{-1}$$

$$\hat{\boldsymbol{\beta}}^{(m)} = \hat{\boldsymbol{\beta}}^{(m-1)} + \left\{ \boldsymbol{X}^{*} \boldsymbol{W} \; \boldsymbol{F}^{-1} (\boldsymbol{\chi} - \boldsymbol{\dot{\boldsymbol{\mu}}}) \right\}_{\boldsymbol{\beta} = \boldsymbol{\hat{\boldsymbol{\beta}}}^{(m-1)}} \left[ (\boldsymbol{X}^{*} \boldsymbol{W} \; \boldsymbol{X})^{-1} \right]_{\boldsymbol{\beta} = \boldsymbol{\hat{\boldsymbol{\beta}}}^{(m-1)}}, (26)$$

Equation (28) is same as

•

$$\left[ (\mathbf{X} \cdot \mathbf{W} \mathbf{X}) \right]_{\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}^{(m-1)}} \hat{\boldsymbol{\beta}}^{(m)} = \left[ (\mathbf{X} \cdot \mathbf{W} \mathbf{X}) \right]_{\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}^{(m-1)}} \hat{\boldsymbol{\beta}}^{(m-1)}$$
$$+ \left\{ \mathbf{X} \cdot \mathbf{W} \mathbf{F}^{-1} (\boldsymbol{\gamma} - \boldsymbol{\mu}) \right\}_{\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}^{(m-1)}}$$

= X'W  $\hat{\underline{z}}^{(m-\pm)}$ ;

where

•

$$\underline{z} = X\underline{\partial} + \left\{ F^{-1}(\underline{y} - \underline{\mu}) \right\}$$

i.e. 
$$\underline{z} = \underline{T} + \left\{ F^{-1}(\underline{y} - \underline{\mu}) \right\}$$

Thus m<sup>th</sup> approximation of estimate of  $\beta$  can be obtained from the equation

$$(X' W X)\hat{\beta} = X' W \hat{z}, \qquad (29)$$

where the matrix W is obtained from current estimate of  $\beta$ . <u>Remark</u>: From equations (7) and (29) it is observed that for generalised linear model with density (or mass) function as in equation (3.2-9), solution of the maximum likelihood equations is equivalent to an iterative weighted least squares procedure with a weight function

$$\mathbf{W} = \mathbf{V}^{-1}\mathbf{F}^2, \tag{30}$$

÷μ

1

and a modified dependent variate

$$\underline{Z} = \hat{\underline{T}} + (\underline{Y} - \underline{\mu}) F^{-1}; \qquad (31)$$

where

(i)  $\mu = E(\underline{Y})$ (ii)  $\Psi = diag(v_{11}, v_{22}, ..., v_{nn})$ , and (iii)  $F = diag(f_{11}, f_{22}, ..., f_{nn})$ , with

 $\begin{cases} f_{ii} = (d\mu_i/dT_i) \\ v_{ii} = var(Y_i)/\alpha(\phi) \end{cases}, \text{ for } i=1,2,\ldots,n.$ 

In fitting generalised linear model, we have used Fisher's scoring method as a numerical technique to obtain parameter estimates. This method was introduced by Fisher(1935), in the appendix of a paper by Bliss (1935). Green(1984), Jorgenson (1984), Finney (1987) and McCuliagh & Nelder (1989) used the same method to fit generalised linear model. Below we give an algorithm of the same method to compute the maximum likelihood estimates of  $\beta$ .

```
3.4.3 : <u>Algorithm to obtain maximum likelihood estimates of </u>B :
While fitting generalised linear model, approximate values
```

of maximum likelihood estimates of  $\beta$  can be obtained by using weighted least square method as a numerical technique. As stated earlier, since the new dependent variable Z and W both depend on the fitted values, the procedure is iterative. Hence estimates should be obtained iteratively by using N-R-method or Fisher's method of scoring. To obtain estimates of  $\beta$ , equation (29) can be used. This method has the following steps.

- (I) Write the incidence matrix X ;
- (II) fix a small positive number  $\varepsilon$  (say, ) to get estimate of  $\beta$  with desired accuracy ;
- (III) find expressions for  $(d\mu_i/dT_i)$  and  $V_{ii}$ ;
- (IV) take averages of different samples in the data as initial estimates  $\hat{\mu}_i^{(O)}$  of  $\mu_i$ , (i=1,2,...,n);
- (V) obtain  $\hat{T}_{i}^{(o)} = m(\hat{\mu}_{i}^{(o)})$ ;
- (VI) compute

$$\hat{W}_{ii}^{(o)} = \begin{bmatrix} 1 \\ ---, (d\mu_i/dT_i)^2 \\ V_{ii} \end{bmatrix}_{\mu_i = \hat{\mu}_i^{(o)}}$$

and

$$\hat{z}_{i}^{(o)} = \hat{T}_{i}^{(o)} + \begin{bmatrix} (y_{i} - \mu_{i}) \\ ----- \\ V_{i} \end{bmatrix} \mu_{i} = \hat{\mu}_{i}^{(o)}$$

(VII) compute  $\hat{\beta}^{(o)}$  as

$$\hat{\beta}^{(0)} = (X, W_{(0)}X)^{-1}X, W_{(0)}^{(0)};$$

(VIII) obtain  $\hat{\underline{T}}^{(\pm)} = \mathbf{X} \hat{\boldsymbol{\beta}}^{(0)}$ , and  $\hat{\boldsymbol{\mu}}_{i}^{(\pm)} = m^{-\pm}(\hat{\underline{T}}_{i}^{(\pm)});$ 

(IX) repeat steps similar to the steps (VI) to (VIII) until

:: .::  $\beta_{j}^{(1)} - \beta_{j}^{(1-1)} < c, \text{ for all } j = 0, 1, \dots, k; \quad (32)$ 

(X)  $\beta_j^{(p)}$ , (j = 0,1,...,k) are the final estimates of  $\beta_j$  if conditions (32) are satisfied for  $\beta_j^{(p)}$ , but are not satisfied for  $\beta_j^{(p-1)}$ , for atleast one value of j.

<u>Remarks</u> : (i) When any of the  $\hat{\mu}_i^{(o)}$  takes extreme value, some adjustment should be done by adding or subtracting proper positive number, so that the corrected initial estimate will not be an extreme value.

(ii) Another important fact to be noted here is about existence of inverse of the matrix (X'W X). Wedderburn (1976) has proved that  $(X'W X)^{-1}$  exists for log concave link functions.

After fitting generalised linear model to the data, one may be interestd in testing different hypotheses about the model parameters  $\beta$ . For this purpose it is necessary to obtain sampling distribution of maximum likelihood estimates of  $\beta$ . 3.4.4 : <u>Sampling dustribution of maximum likelihood estimates</u> :

Suppose (X'W X) is non-singular, so that there are unique maximum likelihood estimates of the parameters  $\beta$  and are close to the true value  $\beta$ . Suppose I is the information matrix. Then using Taylor's first order approximation about  $\hat{\beta}$  for  $\underline{T}^{*}(\beta)$ , we get,

$$\underline{\mathbf{T}}^{*}(\hat{\boldsymbol{\theta}}) = \underline{\mathbf{T}}^{*}(\hat{\boldsymbol{\theta}}) + \mathbf{H}(\hat{\boldsymbol{\theta}})(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}),$$
 (33)

For simplicity, let

$$\underline{U} = \partial l / \partial \beta, \qquad (34)$$

3

: ...

and H is the matrix whose  $(j,i)^{th}$  component is  $(\partial^2 l/\partial \beta_j \partial \beta_l)$ . Then result similar to (3.2-25) gives

$$I = E(U U'),$$
  
= -E(H). (35)

**,**1

ş

As H is asymptotically equal to its expected value and since  $\underline{T}^*(\hat{\beta}) = 0$ , for large samples equation (33) becomes,

which implies,

$$(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) = \mathbf{I}^{-1}\{\underline{\mathbf{T}}^{*}(\boldsymbol{\beta})\}, \qquad (36)$$

If I is the matrix of constants, taking expectation of the equation (36), we have,

$$E(\hat{B} - B) = Q, \qquad (37)$$

and

$$E[(\hat{\theta} - \theta), (\hat{\theta} - \theta)'] = E\{I^{-1}\underline{T}^{*}(\theta), \underline{T}^{*}, (\theta)I^{-1}\}.$$
 (38)

Hence by using equation (36) and since it is assumed that I is the matrix of constants, equation (38) can be rewritten as

 $E[(\hat{\beta} - \beta), (\hat{\beta} - \beta)] = I^{-1},$  (39)

Thus for large samples,

$$(\hat{\partial} - \partial)$$
 has  $N_{(k+1)}(\underline{O}, \mathbf{I}^{-1})$  distribution.

This shows that for large samples,

$$\hat{\mathcal{B}}$$
 has  $N_{(k+1)}(\mathcal{B}, \mathbf{I}^{-1})$  distribution.

i.e.  $(\hat{\theta} - \theta)'I(\hat{\theta} - \theta)$  has  $\chi^2$  distribution with (k+1) d.f.

## 3.5 : Measures of adequacy of the fitted model :

Once the model is proposed, the next part is to see whether

it fits 'well' to the data; i.e. to check how 'good' fitted model describes the data.

A 'good' model has to balance two requirements.

- 1. The model should be as complex enough to approximate the real world phenomenon it describes.
- 2. The model should be as simple as possible for the reason that, simpler it is, the more comprehensible it is.

Thus if there are two models that give approximately the same degree of agreement with reality, we should prefer the 'simpler' model. A model is simpler, if it contains fewer number of parameters. The 'full' model describes the data in the 'best' possible way, but it does not reduces the data as it has the number of parameters equal to the number of observations. This model describes the data in the best possible way because, it assigns complete variation in values of the response variate to the systematic components. On the other hand, 'null' model is the simplest model containing a single parameter. Thus it considers all the variation between values of Y due to the random component. Hence it describes the data in the 'worst' manner. This shows that model should be intermediate model containing p (1 parameters, and describing the data in sufficientlySome well known techniques used to check the better way. adequacy of the fitted model are as follows.

- 1. usual chi-square statistic for testing goodness of fit,
- 2. Pearson's chi-square statistic  $(X^2)$ ,
- 3. Deviance (D).

3.5.1 : <u>Chi-square statistic for goodness</u> of fit : The usual  $\chi^2$ -test of goodness of fit can be used to test the goodness fit of the model. This test uses the  $\chi^2$  statistic given by,

$$\chi^{2} = \sum_{i} \begin{bmatrix} (0_{i} - E_{i})^{2} \\ - - - \frac{i}{E_{i}} \end{bmatrix}, \quad i = 1, 2, ..., n , \qquad (1)$$

which has chi square distribution with (n-t-1) d.f. Here

- (i) t = number of d.f. lost in pooling;
- (ii)  $0_{i}$  = observed frequency of i<sup>th</sup> class;
- (iii)  $E_i = expected frequency of i<sup>th</sup> class.$

.

3.5.2 : <u>Pearson's  $\chi^2$  statistic</u> : Other important measure of goodness of fit is 'generalised Pearson  $\chi^2$  statistic'. This statistic is obtained by using the formula,

$$X^{2} = \sum_{i} \left\{ \frac{(y_{i} - \hat{\mu}_{i})^{2}}{V_{i} (\hat{\mu}_{i})} \right\}, \qquad (2)$$

•. .

Formulae for generalised Pearson  $X^2$  statistic for distributions in table (3.2) are summarised in table (3.3) given below. These formulae can be obtained easily.

TABLE 3.3

Distribution	Pearson X <sup>2</sup> statistic
Normal	$\sum_{i} (y_i - \hat{\mu}_i)^2$
Gamma	$\sum_{i} (y_i - \hat{\mu}_i)^2 / (\hat{\mu}_i)^2$
Exponential	$\sum_{i} (y_{i} - \hat{\mu}_{i})^{2} / \hat{\mu}_{i}^{2}$
Poisson	$\Sigma (y_i - \hat{\mu}_i)^2 / \hat{\mu}_i$
Grouped binary	$\sum (y_{i} - \hat{\mu}_{i})^{2} / (\hat{\mu}_{i} (1 - \hat{\mu}_{i}))$

3.5.3 : <u>Deviance</u> : To assess the goodness of fit of a generalised linear model, the following statistic namely, 'deviance' (D, say) is suggested by Nelder & Wedderburn (1972), which is based on likelihood ratio statistic. Bishop & others (1975) named the same statistic as  $G^2$  statistic. The full model is useful in defining this test statistic.

Let  $\ell(\mu, \phi; y)$  and  $\ell(y, \phi; y)$  be the maximised log likelihoods corresponding to the model under study and the full model respectively. Then the deviance (D) is given by

$$D(y_{i}, \hat{\mu}) = \phi \{-2[\ell(\hat{\mu}, \phi; y) - \ell(y_{i}, \phi; y)]\}, \quad (3)$$

In generalised linear models deviance plays a role similar to the role of residual sum of squares in classical linear models.

To make meaning of the term deviance more clear, below formulae for deviance corresponding to the distributions in table(3.3) are obtained by assuming that there are n observations on Y.

<u>lllustration</u> i :<u>Normal distribution</u> : Let observations  $(Y_i)$  on the response variate Y are independently Normally distributed random variables with parameters  $(\mu_i, \sigma^2)$ . Then from equation

(3.2-4) we have

$$\ell(\hat{\mu}, \sigma^2; \underline{y}) = (-\pi/2) \ln(2\Pi\sigma^2) - (1/2\sigma^2) \sum_i (y_i - \hat{\mu}_i)^2.$$
(4)

Similarly we obtain the formula for  $l(\underline{\gamma}, \phi; \underline{\gamma})$ . Hence from equation (3), deviance becomes

$$D(y, \hat{\mu}) = \sum_{i} (y_{i} - \hat{\mu}_{i})^{2}$$
 (5)

<u>Remark</u> : As already pointed in chapter 2, the goodness of fit of classical linear model with error components distributed as

normal, is based on residual sum of squares. It can be seen that the deviance (D) in equation (5) is nothing but residual sum of squares (E), given in equation (2.4-2).

<u>Illustration</u> 2 : <u>Gamma distribution</u> : Suppose we have n observations  $(Y_i)$  on the response variate Y having gamma distribution with p.d.f. as in equation (3.2-11). Then equation (3.2-13) gives

$$l(\hat{\mu},\nu;\underline{\nu}) = \nu \sum_{i} \left[ -y_{i}/\hat{\mu}_{i} - \ln(\hat{\mu}_{i}) \right] - n\ln(\Gamma(\nu) + n\nu\ln(\nu) + (\nu-1)\sum_{i} \ln(y_{i})$$
(6)

Hence after getting  $l(y, \phi; y)$ , from equation (3) we have,

$$D(\underline{y}; \hat{\mu}) = 2\{\sum_{i} [-l_{\Pi}(y_{i}/\hat{\mu}_{i}) + (y_{i}-\hat{\mu}_{i})/\hat{\mu}_{i}]\}$$
(7)

Nelder & Wedderburn (1972) have shown that the second term in equation (7) is identically equal to zero. Proof is made available in Appendix-1. Thus equation (7) is equivalent to  $D(\underline{y}; \hat{\mu}) = -2 \sum_{i} ln(y_{i} / \hat{\mu}_{i}).$  (8)

<u>Illustration 3 : Exponential distribution :</u>

Suppose Y<sub>i</sub> (i=1,2,...,n) are n independent observations on Y having exponential distribution with mean  $\mu$ . Then the log likelihood of Y can be written as

$$l(\mu; \underline{v}) = -\sum_{i} [ln(\mu_{i}) + (y_{i}/\mu_{i})].$$
(9)

Hence from equation (3), we obtain

$$D(\underline{y}; \hat{\mu}) = -2\{\sum_{i} [l_{n}(y_{i}/\hat{\mu}_{i}) + (y_{i}-\hat{\mu}_{i})/\hat{\mu}_{i}]\}, \quad (10)$$

According to the theorem proved in Appendix-1, it can be seen that the second term in equation (10) is zero. Hence, equation (10) becomes

$$D(y;\hat{\mu}) = -2 \{ \sum_{i} ln(y_{i}/\hat{\mu}_{i}) \}, \qquad (11)$$

<u>Illustration</u> 4 : <u>Poisson</u> <u>distribution</u> :

If we assume Poisson distribution with parameter  $\mu$ , then from equation (3.2-16) it can be seen that

$$\ell(\hat{\mu}, \underline{y}) = \sum_{i} \left[ -\hat{\mu}_{i} + y_{i} \ln(\hat{\mu}_{i}) - \ln(y_{i}!) \right], \qquad (12)$$

This implies,

$$D(y_{i}\hat{\mu}) = 2\{\sum_{i} [y_{i}^{l} n(y_{i}^{j}/\hat{\mu}_{i}^{l}) + (y_{i}^{l}-\hat{\mu}_{i}^{l})]\}, \qquad (13)$$

As per theorem given in Appendix-1, it can be seen that the second term in equation (13) is identically equal to zero. Thus equation (13) is same as

$$D(y_{i}\hat{\mu}) = 2 \sum_{i} y_{i} ln(y_{i}/\hat{\mu}_{i}).$$
 (14)

<u>Illustration 5: Binary distribution for grouped data :</u>

Considering Binomial distribution with parameters  $(m_i^{\sharp}, \mu_i)$  for  $Y_i^{\sharp}=m_i^{\sharp}Y_i$ , where  $Y_i$  is the i observation on Y. Then assuming independence, we have,

$$\mathcal{E}(\hat{\mu}, \underline{m}, \underline{v}) = \sum m_{i}^{*} \{y_{i} \ln(\hat{\mu}_{i}) + (1 - y_{i}) \ln(1 - \hat{\mu}_{i})\} + \sigma(\underline{m}, \underline{v}), \quad (15)$$

where  $c(\underline{m}^{*},\underline{v})$  is a function free from  $\hat{\mu}_{i}$ . Thus from equation (3), deviance is

$$D(y,\hat{\mu}) = 2\sum m_{i}^{*} [y_{i} \ln(y_{i}/\hat{\mu}_{i}) + (1-y_{i}) \ln(1-y_{i})/(1-\hat{\mu}_{i})], \quad (16)$$

Deviances discussed in this section, for different distributions, are summarised in table  $3.2_{\rm H}$ .



Distribution	Veviance
Normal	$\sum_{i} (x_i - \mu_i)^2$
Gamma	-2 Σ[l'n(y, /μ]]
Exponential	$-2 \sum_{i} \{ \ln(y_i / \hat{\mu}_i) \}$
Poisson	2 $\Sigma \{ [y_i \ln(y_i / \mu_i)] \}$
Grouped binäry	$2\Sigma_{i}m_{i}^{*} \{ \{y_{i}, in(y_{i}/\hat{\mu}_{i})\} + (1-y_{i})in\{(1-y_{i})/(1-\hat{\mu}_{i})\} \}$

3.5.4 : Advantages-ofy the deviance (D) statistic :

(1) E statistic 'is\_appropriate for maximum likelihood estimates.

Explanation : From equation (1), it is clear that since  $\ell(\underline{\mu}, \phi, \underline{\nu})$  is maximised log likelihood function for intermediate model, and since D is non negative, maximum likelihood estimates give the minimum value of D. Thus D is appropriate for maximum likelihood estimates.

(11) Conditional break down of D is possible.

<u>Explanation</u>: When there are two models, model (1) and model (2), they are said to be nested if one of them (say e.g. model (2)) contains only a subset of terms contained in other model (i.e. model (1)). For nested models conditional break down of D is possible. In simplest conditional break down of D statistic D(2) corresponding to model (2) can be broken down into two parts;

 (i) a measure of distance of the estimates of parameters in the model (2) from those obtained under model (1);
 (ii) a D statistic D(1) for model (1).

**70** %

For this, rewrite equation (1) as,

$$D(\underline{y}:\hat{\mu}) = \phi \{-2([\ell(\hat{\mu}_{2},\phi;\underline{y}) - \ell(\hat{\mu}_{1},\phi;\underline{y})] + [\ell(\hat{\mu}_{1},\phi;\underline{y}) - \ell(\underline{y},\phi;\underline{y})])\}$$
$$= D[(2)[(1)] + D(1), \qquad (17)$$

where  $\ell(\underline{\mu}_i, \phi; \underline{\gamma})$  denotes the maximised log likelihood under model (1) and D[(2)[(1)] is the conditional D statistic for model (2) given model (1). If model (1) is a full model then D[(2)[(1)] is equal to D(2). Such type of conditional break down does not exist for  $\chi^2$  statistic given in equation (16).

(iii) Structural breakdown of D is possible. The structural break down which is possible with D, is not possible with  $\chi^2$  statistic in equation (16).

To use the deviance(D) as a test statistic for testing goodness of fit of the fitted model, the distribution of 'D' must be known. In classical linear models, as normal distribution is assumed for error components, exact distribution of the deviance can be computed easily. But when we depart from normal distribution and from linearity of different effects, generally exact distribution of the deviance is not obtainable. In some situations like exponential distriution. exact sampling distribution can be achieved. When exact distribution can not be obtained, a chi-square distribution is a better approximation for the difference in deviances.

Sampling distribution of the difference between deviances can be obtained by using that for maximum likelihood estimates of the model parameters  $\beta$ . We have already shown that, for large samples,

 $(\hat{\beta} - \beta)' I(\hat{\beta} - \beta)$  has  $\chi^2$  distribution with (k+1) d.f.

Now, on expanding log likelihood function  $\ell(\beta; \underline{y})$  about the maximum likelihood estimates  $\hat{\beta}$  of  $\beta$ , we get,

 $\ell(\beta; \chi) = \ell(\hat{\beta}; \chi) + (\beta - \hat{\beta}) \cdot \underline{U}(\hat{\beta}) + (1/2)(\beta - \hat{\beta}) \cdot \underline{H}(\hat{\beta})(\beta - \hat{\beta}),$  (18) where  $\underline{U}(\hat{\beta})$  and  $\underline{H}(\hat{\beta})$  are values of  $\underline{U}$  and  $\underline{H}$  evaluated at  $\hat{\beta}$ .

Since  $\hat{\beta}$  are maximum likelihood estimates of  $\beta$ , these are the solutions of the equations  $\underline{U}(\beta) = \underline{0}$ . Hence,

$$\underline{U}(\underline{\beta}) = \underline{0}.$$
 (19)

Using equation (19) and approximating  $H(\hat{\beta})$  by E(H), equation (18) becomes

$$2[\ell(\hat{\boldsymbol{\beta}};\boldsymbol{\chi}) - \ell(\boldsymbol{\beta};\boldsymbol{\chi})] = (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^{*}\mathbf{I}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}).$$

Therefore, from the distribution of  $\hat{\beta}$ , it is clear that for large n, the distribution of  $2[\ell(\hat{\beta};\underline{\gamma}) - \ell(\hat{\beta};\underline{\gamma})]$  has  $\chi^2$  distribution with (k+1) d.f. Thus the distribution of deviance is  $\chi^2$  distribution with (k+1) d.f.

Though this chi-square approximation is not adequate in all the cases, a better approximation is yet to be suggested. The table of deviance differences has its importance to select the terms showing significant effect on value of the response variate.

#### 3.6 : Analysis of deviance

·. :

For orthogonal data with normal errors, analysis of variance (ANOVA) is a very useful statistical tool for separating the effects due to systematic components from those due to random components. Nelder & Wedderburn (1972) suggested a generalisation needed, so that it is applicable for analysing generalised linear models. While making generalisation, two problems should be considered. First is, terms in the model are

generally non orthogonal and the other is, for non normal distributions, different sums of squares (SS) are not measuring properly the effects due to different components.

look towards Before going to the generalisation, we the usual ANOVA table from different angle. ANOVA can be considered as the first difference of the measures of descripancy for a sequence of models each including one term more than the previous model; e.g. in factorial model with two factors A and B (say), we have ANOVA with three terms namely main effects A and B, and the interaction effect A.B. SS for these factors are the first differences of the SSEs corresponding to the sequence of four models 1, A, A+B and A+B+A.B respectively.Note that measure of descripancy for model 1 is just TSS, and that for full model A+B+A.B is O. It is also clear that for the full model d.f.are. . ( • , . · equal to 0.

Now the generalisation is very much clear. Nelder & have used D statistic as a measure of Wedderburn (1972) descripancy for a given sequence of nested generalised linear models, and taking the first differences formed the analysis of Since in generalised linear model, deviance (ANODEV) table. mostly the data is non orthogonal, the interpretation is a bit Each number in the ANODEV table represents complicated. variation due to that after eliminating effects of the terms above it and ignoring effects of the terms below it.

In model fitting problem interpretation of the fitted model is also a very important part. In the next section this part is introduced in brief.

# 3.7 : Interpretation of the fitted model :

ł

Interpretation of the fitted model is also a part and parcel of the model fitting procedure. In other words a model should be

such that the practical conclusions can be drawn from estimated parameters of the model. Thus the question is 'what do the estimated coefficients of the model tells us about research qestion that motivated the study ?' Many times coefficients associated with the explanatory variates are of interest. Very rarely intercept  $\beta_0$  is of interest. Estimated coefficients of the stimulus variates represents rate of change in the value of a function of the response variate. Corresponding to per unit change in the value of stimulus variate. Thus interpretation of data involves two parts.

- (1) Determining functional relationship between the response variate and the stimulus variates  $X = (X_1, X_2, ..., X_k)$ ,
- (2) Defining appropriately, unit of change for the independent variables.

As we have seen earlier, the functional relationship between Y and  $\underline{X}$  is known as link function. Along with the likelihood function, this function is also required for fitting generalised linear model to the data.

Another important segment of the model building is model checking. When some proposed model is fitted to the data, it is necessary to check whether the fitted model is appropriate. This checking is needed, because some times model fits well to the data, though the assumptions made while fitting are invalid. In the next section, discussion is made on model checking for generalised linear models. Model checking includes checking goodness of fit of the fitted model and checking for validity of the assumptions made while formulating the model.

#### 3.8 : Model checking for generalised linear models :

The main problem in data analysis through fitting a model is to select a proper class of the models, so that conclusions drawn

from the analysis are not far away from the truth. Α statistician can choose the model class carefully by taking into account the type and structure of the data. Many times it may happen that, though the class is selected carefully, the data themselves indicate that the model selected is unappropriate. This situation occurs because of two reasons. First, the complete data indicates departure from the fitted model, and secondly, there may be few data points away from the rest which are known as outliers. Below we club the various methods of model checking into three groups, namely,

(I) visual display;

(II) tests of deviations in particular direction;

(III) searching for influential points (outliers).

The methods in the group (I) are similar to the methods of drawing residual plots as described in section (2.5). First we discuss these methods.

As in classical linear models, here also the raw material for model checking is, fitted values of the linear predictor  $\underline{T}$ , new adjusted dependent variable Z and the projection matrix H. The residuals can also be used as a raw material for model checking. The various types of residuals in generalised linear model are discussed below.

3.8.1 : <u>Regiduals</u> : As stated in section (2.5), the residual analysis is essential to test adequacy of the fitted model. Residuals are useful to examine thoroughly the appropriateness of the fitted model. These are also useful to check whether the outliers are present. Hence in the next <u>Part</u> we study the residuals related to generalised linear models.

In case of generalised linear models, the generalisation of residuals is to be done so that it is useful for non Normal

distributions also. Below we discuss three different residuals, suitable for generalised linear models.

<u>Definition-4</u> : <u>Pearson residual</u> :- It is defined as the signed square-root of the component of the Pearson's chi square statistic of goodness of fit. Thus it is given by,

$$r_{p} = (y - \hat{\mu}) / [V(\hat{\mu})]^{-1},$$
 (1)

An important disadvantage of Pearson residual is that its distribution for non-Normal distribution is skewed. So the properties of Normal distribution turns out to be invalid. Hence this residual is not much useful. The residual which is more appropriate is Anscombe residual. This residual was introduced by Anscombe (1953).

<u>Definition-5</u> :<u>Anscombe residual</u>:- Anscombe defined the residual by replacing y in Pearson residual by the function A(y) and  $\hat{\mu}$  by  $A(\hat{\mu})$ . This function A(.) is selected such that the distribution of A(y) is very close to the Normal distribution. Wedderburn proved that, for generalised linear model, the function A(.) is given by,

$$A(.) = \int_{-\frac{1}{\sqrt{\mu}}}^{\frac{d\mu}{4}} \frac{d\mu}{4} \, . \qquad (2)$$

This replacement normalises the probability function, but to stabilize the variance it is necessary to scale it by estimated standard deviation of A(Y). Since it is not easy to compute exact variance of A(Y), we take its first order approximated value.

$$[A'(\hat{\mu})] \{V(\hat{\mu})\}^{(1/2)}, \qquad (3)$$

Thus Anscombe residual becomes,

Below we obtain expression for 'Anscombe residual' corresponding to Poisson distribution.

Illustration 1. Poisson distribution :

Let the single observation Y be having Poisson distribution with mean  $\mu$ . From table (3.2) we have,

$$V(\mu) = \mu_{\mu}$$

which gives,

.

$$A(\mu) = \int_{-\frac{1}{\mu}}^{\frac{d\mu}{1 \times 2}} \mu^{\frac{d\mu}{2 \times 2}},$$

$$= (3/2)\mu^{\frac{d\mu}{2 \times 2}},$$
(5)

and

$$A^{*}(\mu)\{V(\mu)\}^{(1/2)} = \mu^{(1/6)}, \qquad (6)$$

Using equations (5) and (6) in (4), the Anscombe residual for Poisson distribution becomes,

$$r_{A} = (3/2)(y-\hat{\mu})/\{\hat{\mu}^{(1/d)}\}$$
 (7)

Similarly, we can easily obtain formulae for Anscombe residual corresponding to normal, gamma and exponential distributions. To obtain an expression for Anscombe residual corresponding to binomial responses is complicated. Cox & Snell (1968) obtained this expression.

When deviance is used as a measure of goodness of fit, then it is better to use deviance residual. <u>Definition-6</u>: <u>Deviance residual</u> :- It is defined as the product of positive square root of a quantity (d) contributed by each unit for deviance and sign of the difference between y and  $\mu$ . Therefore deviance residual (r<sub>n</sub>) is given by,

$$r_{\rm D} = {\rm sign}(y_i - \mu_i)[d_i]^{(1/2)}.$$
 (8)

The three residuals for some well known distributions from one parameter natural exponential family are tabulated below in table (3.5). These formulae are quite clear from expressions for  $\mu_{i}$ ,  $V_{ii}$ , equations (1) to (4) and equation (8).

Distribution	Pearson Residual (r ) P	Anscombe Residual (r_)	Deviance Residual (r <sub>D</sub> )
Normal	(y <sub>i</sub> - µ <sub>i</sub> )	(y <sub>i</sub> - $\hat{\mu}_i$ )	{Sign( $y_i - \hat{\mu}_i$ )}. $ y_i - \hat{\mu}_i $
Gamma	(y <sub>i</sub> - μ <sub>i</sub> )/ μ <sub>i</sub>	$\frac{3(y_{i}^{1/2} - \hat{\mu}_{i}^{1/2})}{\hat{\mu}_{i}^{1/2}}$	{Sign(y <sub>i</sub> - $\hat{\mu}_i$ )}. {-2[ln(y <sub>i</sub> / $\hat{\mu}_i$ )] + (y <sub>i</sub> / $\hat{\mu}_i$ )-1} <sup>1/2</sup>
Exponential	(y <sub>i</sub> - μ <sub>i</sub> )/ μ <sub>i</sub>	$\frac{3(y_{i}^{1/2} - \hat{\mu}_{i}^{1/2})}{\hat{\mu}_{i}^{1/2}}$	{Sign(y <sub>i</sub> - $\hat{\mu}_{i}$ )}. {-2[ln(y <sub>i</sub> / $\hat{\mu}_{i}$ )] - (y <sub>i</sub> / $\hat{\mu}_{i}$ )+i} <sup>4/2</sup>
Poisson	$(y_i - \hat{\mu}_i) / \hat{\mu}_i^{1/2}$	$\frac{3(y_{i}^{2/3} - \hat{\mu}_{i}^{2/3})}{2(\hat{\mu}_{i}^{4/6})}$	{Sign( $y_i - \mu_i$ )}. {2[ $y_i \ln(y_i / \hat{\mu}_i)$ ] - $y_i + \mu_i$ } <sup>1/2</sup>

TABLE	3.5

Though the formulae for Anscombe residual and deviance

residual are looking different for non-Normal distributions, 'for given values of y and  $\mu_{i}$ , these two residuals are having VELA similar values. By using any of these residuals one can carry out the residual analysis. These residuals are in terms of  $\mu$ . The replacement of  $\mu$  by  $\mu$  requires a standardisation factor. The standardisation of Pearson's residual is discussed by McCullagh & Nelder (1983) and that of deviance residual by Cox & Snell (1960). If absolute value of the residual for some data point is - greater than '2', the corresponding data point should be checked for outlier's test. We will not discuss here the methods for group (III) due to their vastness. Now we discuss the methods in group (11) briefly.

(II) Tests of deviations in particular directions :

After fitting the generalised linear model, the most important question comes in the mind is that can the value of deviance function decrease significantly by

- (a) including an extra stimulus variate,
- (b) changing scale of the stimulus variates,
- (c) changing link function in a particular direction,
- (d) changing the variance function,

with the help of available information about the stimulus variates, link function and variance function.

(a) <u>Selection of covariates</u> : Generally the experiment contains information on large number of stimulus variates, but a model should contain less number of parameters and should fit sufficiently well to the data. So it is necessary to select a set of useful covariates. Then sequence of nested generalised linear models is fitted and using the deviance function it is decided, which covariate shows significant effect on the response variate. Covariates showing significant effect will be included

in the model.

(b) Checking scale of stimulus variate : Method of checking scale of the stimulus variate has following steps.

(1) Suppose X and B denotes respectively the covariate whose scale is to be checked and coefficient of the covariate.
 (1) Replace X by

$$g^{*}(\Theta;X) = \begin{cases} X^{\Theta}; & \text{for } \Theta \neq 0, \\ \lambda n'(X); & \text{otherwise.} \end{cases}$$

ţ

9)

(III) Fix  $\varepsilon > 0$ , a small positive number, to obtain estimate of  $\theta$  with desired accuracy.

- (IV) To start calculations take initial value ( $\theta^{(0)}$ ) of  $\theta$  as unity.
- (V) As Taylor series expansion for  $g^{*}(\Theta, X)$ , with first order approximation gives

$$g^{*}(\Theta, X) = g^{*}(\Theta^{(O)}, X) + (\Theta - \Theta^{(O)}) [\partial g^{*} / \partial \Theta]_{\Theta = \Theta^{(O)}},$$

replace  $\beta g^{*}(\Theta, X)$  by two linear terms  $\beta U + \gamma V$ ; where,

$$U^{(O)} = g^{*}(\theta^{(O)}; X),$$
  
$$V^{(O)} = \frac{\partial g^{*}(\theta; X)}{\partial \theta} = \frac{\partial g^{*}(\theta; X)}{\partial \theta}$$

and

$$\gamma^{(0)} = \beta(\theta - \theta^{(0)}),$$

(VI) Fit generalised linear model with U and V as covariates.

(VII) Compute  $\theta^{(1)}$ , an improved estimate of  $\theta$  by

$$\theta^{(1)} = \theta^{(0)} + \gamma^{(0)} / \beta^{(0)},$$

(VIII) Repeat steps similar to the steps (V) to (VII) until

.

$$\theta^{(l)} - \theta^{(l-s)} < \varepsilon$$
.

(IX) Final estimate of  $\theta$  is  $\theta^{(1)}$ .

<u>Remarks</u> : (1) If the initial estimate  $\theta^{(0)}$  is far away from the true value of  $\theta$ , convergence of the process is not sure.

(2) Though this process is very useful, it is not a good technique to include more non linear parameters in the model, when other covariates are highly correlated. This is so because, generally estimates of the non linear parameters have large sampling error and are highly correlated with each other and with corresponding linear parameters.

(c) <u>Checking Link Functions</u> : While fitting generalised linear model, link functions are assumed to be known. Instead of this, it is useful to assume that, link functions come from a class of link functions and particular value of one or more parameters describes elements of that class. Most of the times a class of one parameter link functions is either taken as

$$T = \begin{cases} \mu^{\Theta}; & \text{if } \Theta \neq 0, \\ \ln(\mu); & \text{otherwise}; \end{cases}$$
(10)

or

## $T = \{ [\mu^{\theta} - 1] / \theta \}.$

Pregibon (1980) proposed linearising technique to get optimum estimate of  $\theta$ . The procedure described by him has

following steps.

- (1) Suppose  $T = m^*(\Theta, \mu)$  is the link function, as defined in equation (10).
- (11) Fix  $\varepsilon > 0$ , a small number, to estimate  $\theta$  with desired accuracy.
- (III) To start calculations take initial value  $(\theta^{(0)})$  of  $\theta$  as unity.
- (IV) As Taylor series expansion for  $m^{\frac{1}{2}}(\Theta,\mu)$ , with first order approximation gives

$$\mathfrak{m}^{*}(\Theta, \mathfrak{U}) = \mathfrak{m}^{*}(\Theta^{(O)}, \mathfrak{U}) + (\Theta - \Theta^{(O)})[\cdot \partial \mathfrak{m}^{*}/\partial \Theta]_{\Theta = \Theta^{(O)}},$$

$$m^{*}(\theta,\mu) = \underline{T}^{(0)} + (\theta - \theta^{(0)}) \mu^{\theta^{(0)}} \lim_{\mu \to 0} (\mu),$$
 (11)

Equation (11) can be rewritten as,

$$\underline{T}^{(o)} = \underline{m}^*(\theta, \underline{\mu}) - (\theta - \theta^{(o)}) \underline{\mu}^{\theta^{(o)}} \underline{l}_{\underline{n}}(\underline{\mu}),$$

which is equivalent to,

$$\underline{T}^{(0)} = \sum_{j} X_{i,j} \beta_{j} + \gamma^{(0)} X_{i(k+1)}, \qquad (12)$$

where,

$$X_{i,i,k+1} = \mu_{i}^{\Theta} l_{n}(\mu_{i}^{2}).$$
 (13)

: .

:

- (V) Fit generalised linear model with linear predictor  $\underline{T}^{(O)}$  as given in equation (12).
- (VI) Also fit generalised linear model with linear predictor  $\frac{T}{i} = \sum_{j} X_{ij} \beta_{j}.$

(VII) If difference between deviances of the two fitted

models at steps (V) and (VI) is significant, conclusion is  $\theta^{(O)}$  is appropriate value of  $\theta$ ,

(VIII) Take different values of  $\theta$  in place of  $\theta^{(0)}$  and repeat the steps (IV) to (VII). The value of  $\theta$  for which deviance is minimum, is the maximum likelihood estimate of  $\theta$ .

(d) <u>Checking Variance Function</u>: Neider & Pregibon (1987) suggested a method of comparing variance functions for continuous data by using the idea of 'extended quasi likelihood function'. For the distributions discussed in the table (3.2), it can be shown that the log likelihood (l) is close to an extended quasi likelihood discussed in chapter 5. This fact can be used for checking variance function.

## 3.9 Method of obtaining robust estimates :

As discussed in section (2.3), least absolute deviations approach to estimate the parameters was introduced by Boscovich in the year 1757, about 40 years back to the introduction of least square approach due to Gauss in 1797. Hence it 18 naturally quite intersesting to see whether least absolute deviations approach can be used in generalised linear models, instead of usual weighted least square approach, to estimate the model parameters. Morgenthaler(1992) has explained haw least absolute deviations principle can be used to find robust estimates of the model parameters in case of generalised linear models.

During the last decade or more, several statisticians worked on robust estimation of the model parameters for generalised linear models. Some of them are Pregibon(1982), Stephanski & others(1986), and Kunsch & others(1989). All the above statisticians cocentrated particularly on logistic model (model

is discussed in chapter 4). Morgenthaler(1992) described robust estimation in case of generalised linear models. Below the method of obtaining robust estimates is given explicitly.

Method of obtaining robust estimates :

Since least absolute deviations resists the gross error significantly, it is of interest to study haw least absolute deviations principle can be used to obtain robust estimates. Suppose  $Y_i$  (i=1,2,...,n) are n independent responses with

$$E(Y_{i}) = \mu_{i},$$

$$Var(Y_{i}) = V_{ii}(\mu_{i}).$$
(1)

Then robust estimates for fitting L - norm  $(q \ge 1)$  can be obtained by minimising the quantity,

$$K(\underline{\mu},\underline{Y}) = \sum_{i} \left| \frac{(Y_{i}-\mu_{i})}{(V_{ii}(\mu_{i}))^{-(1/2)}} \right|^{q}$$

Therfore, i<sup>th</sup> component of the gradient corresponding to the quantity  $K(\mu, \underline{Y})$  is

q 
$$[V_{ii}(\mu_i)]^{-(q/2)} |Y_i - \mu_i|^{(q-1)} sgn(Y_i - \mu_i),$$
 (2)  
i=1,2,...,n.

For q = 2, the quantity in (1) is equivalent to the i<sup>th</sup> component of quasi likelihood (quasi likelihood is introduced in chapter 5). Morgenthaler(1992) mentioned that for any other value of q, gradient in (2) non consistent estimates of  $\beta$  when the responses are not symmetrically distributed around their means. To obtain symmetrically distributed responses, the distributional form of the response variate must be known.

Since in generalised linear models, distributional form of the response variate is known, correction factors for gradient given in (2), to obtain consistent estimates of  ${\mathcal B}$  are naturally,

$$C_{i} = E\left\{ |Y_{i} - \mu_{i}|^{(q-1)} \operatorname{sgn}(Y_{i} - \mu_{i}) \right\}, \qquad (3)$$

$$i=1, 2, \dots, n.$$

Hence the i<sup>th</sup> (i=1,2,...,n) component of the corrected gradient is,

$$q [V_{ii}(\mu_i)]^{-(q/2)} [Y_i - \mu_i]^{(q-1)} sgn(Y_i - \mu_i) - C_i].$$

Thus estimating equations to obtain estimates of *B*, corrected for consistency are,

$$q[V_{i}(\mu_{i})]^{-(q/2)}|Y_{i}-\mu_{i}|^{(q-1)}[sgn(Y_{i}-\mu_{i}) - C_{i}](d\mu_{i}/dT_{i})X_{i} = 0,$$
  
i=1,2,...,n. (4)

Hence this method can be used as an alternative method to estimate the model parameters  $\beta$  in case of generalised linear models.

In all the above discussion of generalised linear models, it is assumed, for known distributional form of the responses  $Y_i$ (i=1,2,...,n), variance of  $Y_i$  is a specific function  $V(\mu_i)$  of the mean  $\mu_i$  of  $Y_i$ . Mathematically, the variance of  $Y_i$  (i=1,2,...,n) is given by,

$$Var(Y_i) = \phi V(\mu_i), \qquad (5)$$

Ł

To clearify meaning of the above statement, an illustration is given below.

<u>Illustration</u> 1. <u>Normal distribution</u> :

Suppose Y<sub>1</sub>(i=1,2,...,n) are n independent responses having  $N(\mu_i, \sigma^2)$  distribution. From the table (3.2) we recall that, for  $N(\mu_i, \sigma^2)$  distribution,

 $\phi = \sigma^2$  and  $V(\mu_i) = 1$ .

Thus in this case, the variance of  $Y_i$  (i=1,2,...,n) is a specific function of  $\mu$  up to a muliplicative constant  $\sigma^2$ .

On the other hand, if we assume that responses  $Y_i$ (i=1,2,...,n) are independently distributed  $N(\mu_i,\sigma_i^2)$  variables, then

$$\phi_i = \sigma_i^2$$
 and  $V(\mu_i) = 1$ .

Therefore, in this case  $Var(Y_i)$  is not a specific function of the mean  $\mu_i$  (i=1,2,...,n). This is a case where the factor, related to the dispersion parameter is varying instead of being a constant. In such situations if  $\phi_i$  is unknown, 'generalised linear model' can not be fitted in its original form given by Nelder & Wedderburn (1972). The models which are useful here are 'generalised linear models with varying dispersion.

## 3.10 Generalised linear models with varying dispersion

Smyth (1989) introduced the 'generalised linear model with varying dispersion'. Thus he generalised 'generalised linear models', by including a 'dispersion model' along with the usual 'mean model'. The term 'mean model' which occurs first time, is nothing but the usual 'generalised linear model' given in definition (3). Before Smyth (1989), Nelder & Pregibon (1987) introduced new class of models, namely, 'extended quasi likelihood models'. The model introduced by Smyth (1989) is very much similar to the extended quasi likelihood model. Below we discuss 'generalised linear model with varying dispersion'.

Generalised linear model with varying dispersion' can be defined in the following two parts.

## <u>Definition-7</u> : <u>Mean model</u> :

Let Y be a response variate with p.d.f. (or p.m.f.)  $f(.;\Theta^{T})$  which belongs to the one parameter natural exponential family.

Suppose Y<sub>i</sub> (i=1,2,...,n) are n independent observations on the response variate Y such that  $E(Y_i) = \mu_i$  and  $Var(Y_i) = \phi_i V_{ii}(\mu_i)$ . Let  $\underline{x_1}, \underline{x_2}, \ldots, \underline{x_k}$  be the vectors of known values of the covariates X<sub>i</sub> (j=1,2,...,k). Suppose

$$T_{i} = \beta_{0} + \sum_{j} x_{ij} \beta_{j}$$
, for i=1,2,...,n, (1)

where

 $\beta = (\beta_0, \beta_1, \dots, \beta_k)'$  is the vector unknown model parameters and

 $\underline{T} = (T_1, T_2, \dots, T_n)'$  is the linear predictor.

Then a class of models of the form,

$$\underline{\mathbf{T}} = \mathbf{m}(\mathbf{E}(\underline{\mathbf{Y}})); \tag{2}$$

where m(.) is a strictly monotoic differentiable function; is called as the mean model.

For the 'dispersion model', deviance residuals (d<sub>i</sub>) are generally taken as unobservable responses. Then dispersion model can be defined as below.

## <u>Definition-8</u> : <u>Dispersion</u> model :

Consider the unobservable responses  $d_i$  (i=1,2,...,n) with  $E(d_i) = \phi_i$  and  $Var(d_i) = \theta V_B(\phi_i)$ . Let  $\frac{z^*}{j}$  (j=1,2,...,k<sup>\*</sup>) be the vectors of known values of the covariates  $Z_j^*$  (j=1,2,...,k<sup>\*</sup>). Suppose

$$\zeta_i = \gamma_0 + \sum_j z_{ij}^* \gamma_j$$
, for i=1,2,...,n, (3)

where

 $\chi = (\gamma_0, \gamma_1, \dots, \gamma_k^*)'$  is the vector unknown model parameters and

 $\zeta = (\zeta_1, \zeta_2, \dots, \zeta_n)'$  is the linear predictor for dispersion model.

Then a class of models of the form,

$$\zeta = h(E(\underline{d})); \qquad (4)$$

where h(.) is a strictly monotoic differentiable function; is called as the dispersion model.

<u>Definition-9</u> (Smyth 1989) : <u>Generalised linear model with varying</u> <u>dispersion</u> :

The 'generalised linear model with varying dispersion' can be defined in conjunction with the above definitions of mean model and dispersion model as the class of models of the form

m(
$$\mu$$
)=X $\beta$ , h( $\phi$ )=Z $\gamma$ , Var(Y<sub>i</sub>)= $\phi_i V_{ii}(\mu_i)$  and Var(d<sub>i</sub>)= $\theta_i V_D(\phi_i)$ . (5)  
for i=1.2....n

The following theorem helps in developing the procedure of fitting generalised linear model with varying dispersion.

Theorem 3.4 : Mean and dispersion in the generalised linear model with varying dispersion are orthogonal.

**Proof :** Suppose Y<sub>i</sub> (i=1,2,...,n) are n independent observations on the response variate. Suppose the dispersion parameters  $\phi_i$ (i=1,2,...,n) are not constant for all responses. Then log likelihood of the complete data set is given by,

$$\ell(\underline{\theta}, \underline{\phi}; \underline{y}) = \sum_{i} \left\{ -\frac{1}{\phi_{i}} - \left[ (y_{i} \theta_{i} - g(\theta_{i}) + h(y_{i}) \right] + \xi(\phi_{i}; y_{i}) \right\}$$
(6)

Differentiating equation (6) w.r.t.  $\theta_i$  and  $\phi_i$  (i=1,2,...,n) we have,

$$\partial^{2} \mathcal{U}(\partial \Theta_{i} \partial \phi_{i}) = -(y_{i} - g'(\Theta_{i}))/(\phi_{i}^{2}), \qquad (7')$$

Taking expectation on both the sides of equation (7) and using the fact that  $g'(\Theta_i) = \mu_i$  (for i=1,2,...,n) we get,

 $E(\partial^2 l/(\partial \Theta_i \partial \phi_i)) = 0.$  (8)

The equations (8) above imply that mean and dispersion parameters are orthogonal.

In generalised linear models with varying dispersion since the mean and dispersion parameters are orthogonal, it is possible to estimate the parameters  $\beta$  and  $\gamma$  one at a time. We discuss below the procedure of fitting generalised linear model with varying dispersion.

3.11 Fitting of generalised linear model with varying dispersion: Suppose  $\hat{\mu}^{(r)}$  and  $\hat{\phi}^{(r)}$  (r=0,1,...) denote the  $j^{th}$ approximation to the estimate of  $\mu$  and  $\phi$  respectively. Due to the interlinking of the two models, fitting procedure is alternating as described below.

While fitting the mean model keep  $\phi$  fixed at  $\hat{\phi}$ . Similarly, fit dispersion model by fixing  $\mu$  at  $\hat{\mu}$ . This method has the following steps.

- (1) Decide the 'independent' variates  $Z_{j}^{*}$  (j=1,2,...,k<sup>\*</sup>) for the 'dispersion' model.
- (2) Choose two small positive numbers  $\varepsilon_1$  and  $\varepsilon_2$  according to the desired accuracy for  $\hat{\beta}$  and  $\hat{\gamma}$ .
- (3) To start the calculations take initial estimates of  $\phi$ and  $\mu$  as  $\hat{\phi}^{(0)} = E_{\mu}$  and  $\hat{\mu}^{(0)} = y_{\mu}$
- (4) Fit the mean model as usual by using algorithm discussed in sub section (3.4.3) for fixed  $\hat{\phi}$  (the current estimate of  $\phi$ .
- (5) Compute the deviance residuals  $d_i$  (i=1,2,...,n) from the fitted mean model.
- (6) Fit the dispersion model for  $d_i$  (i=1,2,...,n) by

assuming gamma distribution. While fitting the dispersion model, fix the vector of parameters  $\mu$  at  $\dot{\mu}$  (the current estimate of  $\mu$  for the mean model). (7) Répeat the steps (4) to (6) until,

$$|\gamma_{j}^{(m_{j})} - \gamma_{j}^{(m_{j}-4)}| < \varepsilon_{2}, \text{ for all } j=0,1,\ldots,k^{*}.$$
 (10)

Then final estimates of  $\beta$  and  $\gamma$  are

$$\hat{\beta} = \hat{\beta}^{(m)}$$
 and  $\hat{\gamma} = \hat{\gamma}^{(m)}$ . (11)

By using the above mentioned procedure one can fit the generalised linear model with varying dispersion.

 $|\beta_{i}^{(m)} - \beta_{i}^{(m-1)}| < \epsilon_{i}, \text{ for all } j=0,1,...,k$ 

In many fields like socio-economic field, frequently the data are of discrete type. Hence in the next chapter we discuss the fitting of generalised linear model for various types of discrete data.



11

(9)

90

ેં ડેર