

- boundary sheet
- Goodness of fit
- (log)-likelihood function
- statistic base

- Summary of process & logic of hypothesis tests

 - 1) Specify Model H_0 (corresponding to θ_0)
 $\cdots \cdots \cdots H_0 \cdots \cdots \cdots H_1$
 - 2) Fit H_0 & calculate goodness of fit statistic G_0
 $\cdots M_0 \cdots \cdots \cdots G_0 \cdots \cdots \cdots G_1$
 - 3) Calculate the improvement in fit, usually $G_1 - G_0$ or G_1/G_0
 - 4) Use sampling dist of $G_1 - G_0$ (or whatever from 3) to test hypothesis $G_0 = G_1$ vs $G_0 \neq G_1$.
 - 5) If $G_0 = G_1$ is not rejected \Rightarrow fail to reject H_0 + H_0 is preferred model
 $G_0 \neq G_1$ is rejected \Rightarrow reject H_0 + $H_1 \cdots \cdots \cdots$

→ Sampling distributions

→ For both forms of inference (CI + HT), the sampling distributions are required

$$\begin{matrix} \text{and samp dist} & \text{and samp dist of} \\ \text{of estimator} & \text{goodness of fit statistic} \end{matrix}$$

→ If response \sim normal dist, then samp dist can be determined exactly.
 For others, use asymptotics based on CLT
 (ignoring attention to regularity conditions
 \Rightarrow if $\mathbb{E}(S)$ are II & from exp family \Rightarrow conditions are met)

→ Basic idea is that under appropriate conditions if S is a statistic of interest then approximately

$$\frac{S - E(S)}{\sqrt{V(S)}} \sim N(0,1)$$

or equivalently

$$\frac{(S - E(S))^2}{V(S)} \sim \chi^2_1$$

→ If there is a vector of statistics $S = \begin{pmatrix} S_1 \\ \vdots \\ S_p \end{pmatrix}$
 w/ asymptotic expectation $E(S)$ & var-cov matrix V

$$[S - E(S)]^T V^{-1} [S - E(S)] \sim \chi^2_p$$

(non-singular & inverse exists)

$$\rightarrow \text{Score} = \sum_i y_i \ln p_i + (1 - y_i) \ln (1 - p_i)$$

- Var-Cov matrix of U_i is the information matrix $J(\theta)$ elements $J_{jk} = E[U_j U_k]$

→ If only one parameter β , score stat has asymptotic dist $\frac{U}{\sqrt{J}} \sim N(0,1)$ or equivalently $\frac{U^T}{J} \sim \chi^2(1)$

($E(U) = 0$ & $V(U) = J$)

→ If there is a vector of parameters $\beta = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$ & score vector $U = \begin{bmatrix} U_1 \\ \vdots \\ U_p \end{bmatrix} \sim MVN(0, J)$

$\Rightarrow U^T J^{-1} U \sim \chi^2(p)$ for large samples

→ Example $\rightarrow Y_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$
 ↳ known constant

\rightarrow Log likelihood $\rightarrow l = -\frac{1}{2\sigma^2} \sum (y_i - \mu)^2 - n \ln(1 + e^{\mu})$

\rightarrow Score statistic $\rightarrow U = \frac{\partial l}{\partial \mu} = \frac{1}{\sigma^2} \sum (y_i - \mu) = \frac{n}{\sigma^2} (\bar{Y} - \mu)$ ①

$\rightarrow MLE \rightarrow$ obtained by solving $U = 0 \Rightarrow \hat{\mu} = \bar{Y}$

\rightarrow Expected value & variance $\rightarrow E(U) = \text{cov}(U| \mu) = \frac{1}{\sigma^2} \sum_{i \neq j} E[(Y_i - \mu)(Y_j - \mu)] = 0$

$\rightarrow J = V(U) = \text{cov}(U| \mu) = \frac{1}{\sigma^2} \sum V(Y_i) = \frac{n}{\sigma^2}$

$\Rightarrow \frac{U}{\sqrt{J}} = \text{cov}(U| \mu) = \frac{(\bar{Y} - \mu)}{\sigma/\sqrt{n}} \sim N(0, 1)$
 asymptotically (exactly in this scenario)

Similarly $U^T J^{-1} U = \frac{U^2}{J} = \frac{(\bar{Y} - \mu)^2}{\sigma^2/n} \sim \chi^2(1)$ ✓

\Rightarrow Use this sampling dist of U to make inferences about μ
 e.g.) 95% CI $= \bar{Y} \pm 1.96 \text{SE}_U$

→ Example \rightarrow If $Y \sim Bin(n, p)$

\rightarrow Log-likelihood $\rightarrow l = n \ln(p) + (n-y) \ln(1-p) + \ln(\frac{1}{2})$

\rightarrow Score statistic $\rightarrow U = \frac{\partial l}{\partial p} = \frac{y}{n} - \frac{n-y}{n-p} = \frac{y-np}{np(n-p)}$

\rightarrow Mean & variance $\rightarrow E(U) = np \Rightarrow E(U) = 0$ as expected

$\rightarrow V(U) = n \pi(1-\pi) \Rightarrow V(U) = \frac{1}{np(n-p)} V(Y) = \frac{n}{np(1-p)}$

$\rightarrow \frac{U}{\sqrt{V(U)}} \sim N(0, 1)$

→ Taylor series app
→ used to
for various

- provided π is near π^*

→ for log-likelihood function of single parameter β , the last 3 terms become

$$L(\beta) = l(\beta) + (\beta - \beta_0) u'(\beta) + \frac{1}{2} (\beta - \beta_0)^T U'(\beta)$$

\downarrow ↗ score function evaluated at β

$$U'(\beta) = \frac{\partial^2 l}{\partial \beta^2} \approx \Sigma U_i^2 = -J \quad \text{from earlier formulas} \rightarrow$$

⇒ approximation

$$l(\beta) \approx l(\beta_0) + (\beta - \beta_0) u(\beta_0) - \frac{1}{2} (\beta - \beta_0)^T J(\beta_0)$$

\downarrow ↗ information evaluated at $\beta = \beta_0$

→ for vector β , this becomes

$$l(\beta) = l(\beta_0) + (\beta - \beta_0)^T u(\beta_0) - \frac{1}{2} (\beta - \beta_0)^T J(\beta_0) (\beta - \beta_0)$$

→ for score function of a single parameter β , using first two terms of series, we get

$$u(\beta) = u(\beta_0) + (\beta - \beta_0) \underbrace{u'(\beta_0)}_{U \approx E(u')} = -J$$

\downarrow

$$U = u(\beta_0) - (\beta - \beta_0)^T J(\beta_0)$$

for vector $\Rightarrow u(\beta) = u(\beta_0) - J(\beta_0)(\beta - \beta_0)$ \cancel{A}

5.4 → Sampling distributions for ALEs

→ An equation directly above can be used to find sampling dist. of ALE $b = \hat{\beta}$. By def., b is an estimator which

maximizes $l(\beta)$ & so $u(\beta) = 0$ (first derivative)

⇒ $u(\beta) = -J(\beta) (\beta - \beta_0)$ or equivalently

$$(\beta - \beta_0) = J^{-1} u \quad \text{(leads naturally } J^{-1} \text{ is use negative to switch } b \text{ & } \beta)$$

→ if J is treated as constant

$$E(b - \beta_0) = 0 \quad b/c \quad E(u) = 0$$

$$\Rightarrow E(b) = \beta_0 \Rightarrow b \text{ is a consistent estimator of } \beta$$

(+ least asymptotically)

→ The variance matrix for b is

$$E[(b - \beta_0)(b - \beta_0)^T] = E[(J^{-1}u)(J^{-1}u)^T]$$

\downarrow

$$= J^{-1} \underbrace{E(uu^T)}_{= I} J^{-1} \quad \text{constant} \quad \& (J^{-1})^T = J^{-1} \text{ vs symmetric}$$

$$= J^{-1}$$

⇒ Asymptotic sampling dist for b is

$$(b - \beta_0)^T J(J(b))^{-1} (b - \beta_0) =$$

$$(b - \beta_0)^T J(b) (b - \beta_0) \sim \chi^2(p)$$

→ This is the Wald statistic

→ for one parameter → more commonly used form $b - \pi/\beta, J^{-1}$

→ if response \sim normal \Rightarrow exact result

→ Example → Consider model $E(y_i) = \beta_0 + x_i^T \beta, y_i \sim N(\mu_i, \sigma^2)$

⇒ $y_i \sim \mathcal{I}$ identity link function

→ Information matrix $\rightarrow J_{jk} = \sum_{i=1}^n \frac{\partial^2 \ln L(\beta)}{\partial \beta_j \partial \beta_k} =$

\downarrow

$$= \sum_{i=1}^n \frac{x_{ij} x_{ik}}{\sigma^2}$$

$$= \frac{1}{\sigma^2} x^T x$$

→ $Z_i = \underbrace{\sum_{k=1}^p x_{ik} \beta_k}_{b = x^T \beta} \xrightarrow{\text{evaluated at } \beta = \beta^{(m-1)}} (y_i - \mu_i) = y_i$

$$\Rightarrow x^T w \times b^{(n)} = x^T w \cdot z \text{ is equal to}$$

$$\xrightarrow{\text{with } w} \frac{1}{\sigma^2} x^T x b = \frac{1}{\sigma^2} x^T y$$

$$\xrightarrow{\text{with } w} \Rightarrow \text{MLE is } b = (x^T x)^{-1} x^T y$$

\rightarrow writing model as matrix form \rightarrow MLE MLE (exp. $\sigma^2 I$)

\rightarrow expected value & variance of b

$$\rightarrow E(b) = E[(x^T x)^{-1} x^T y]$$

$$= (x^T x)^{-1} x^T E(y)$$

$$= (x^T x)^{-1} x^T \underbrace{\beta}_{\mathbf{z} = x\beta}$$

$$= \beta \Rightarrow \text{unbiased}$$

\rightarrow Variance of b

$$b - \beta = (x^T x)^{-1} x^T y - \beta$$

$$= (x^T x)^{-1} x^T y - (x^T x)^{-1} (x^T \beta)$$

$$= (x^T x)^{-1} x^T (y - x\beta)$$

$$\Rightarrow E[(b - \beta)(b - \beta)^T] = (x^T x)^{-1} x^T E[(y - x\beta)(y - x\beta)^T] x \xrightarrow{(AB)^T = B^T A^T} (x^T x)^{-1}$$

$$= (x^T x)^{-1} x^T V(y) x \xrightarrow{(x^T x)^{-1}}$$

$$= \sigma^2 (x^T x)^{-1}$$

$$= \sigma^2 I$$

\rightarrow In end \rightarrow MLE b is a linear combination $y_i + \text{Nc random}$

$\xrightarrow{\text{Nc normal}} \Rightarrow$ exact sampling dist
 $\xrightarrow{b \sim N(\beta, \sigma^2 I)}$, or equivalently
 $(b - \beta)^T T(b - \beta) \sim \chi^2_{np}$

5.5 \rightarrow Log-likelihood ratio statistic

\rightarrow One way of assessing the adequacy of a model is to compare it with a more general model w.r.t. the number of parameters that can be estimated.

This is called a saturated (maximal or full) model

\rightarrow If there are N obs y_1, \dots, y_N ($\forall i$) potentially different values of the linear component $x_i^T \beta$, then a saturated model can be specified w.r.t. N parameters

\rightarrow If there are replicates ($\exists i$ w.r.t. some entries of X), then N parameters in full model may be $< N$

\rightarrow Setup \rightarrow Let m be the max. # of parameters that can be estimated
 \rightarrow b_{max} denote the parameter vector for the saturated model
 $\rightarrow b_{\text{max}} = \beta_{\text{max}}$

\rightarrow The likelihood & log-likelihood functions evaluated at b_{max} will be larger than any other likelihood function (using same dist & func function)
 $\&$ it provides the most complete description of the data

\rightarrow Statistic \rightarrow wr $L(b/\gamma)$ by the likelihood function for the model of interest
 \rightarrow Then the likelihood ratio

$$\lambda := \frac{L(b_{\text{max}}/\gamma)}{L(b/\gamma)} \Rightarrow L(\lambda) = L(b_{\text{max}}) - L(b/\gamma)$$

provides a way of assessing the goodness of fit for the model

\rightarrow Large values of λ or $\ln(\lambda)$ \Rightarrow poor description for model
 $\&$ interest relate to saturated model

\rightarrow Need sampling dist of λ to determine critical values

$$\rightarrow B + \frac{D}{\lambda} \sim \chi^2 \Rightarrow$$
 we thus statistic instead
 $\lambda \equiv \text{Deviance}$

5.6 \rightarrow Sampling dist of deviance

\rightarrow Definition \rightarrow Deviance / Log-likelihood (ratio) statistic

$$D = 2[L(b_{\text{max}}/\gamma) - L(b/\gamma)]$$

\rightarrow Distribution \rightarrow If b is MLE of β (\rightarrow s.t. $L(b/\gamma) = 0$)

$$\begin{aligned} (\text{first derivative}) \quad L'(\beta) &= L(b) + (P-b)^T U(b) - \frac{1}{2} (P-b)^T J(b) (P-b) \\ (\text{Taylor series approximation}) \quad \downarrow &\quad \downarrow \\ L'(P) - L(b) &= -\frac{1}{2} (P-b)^T J(b) (P-b) \end{aligned}$$

approximately

$$\Rightarrow 2[L(b) - L(P)] = \frac{1}{2} (P-b)^T J(b) (P-b) \sim \chi^2_{k \text{ of parameters}}$$

(where switched parameter for deviance, but estimate by squaring \Rightarrow not important)

\rightarrow Sampling dist of D

$$\rightarrow D \sim 2[L(b_{\text{max}}) - L(b)]$$

$$\downarrow = 2[L(b_{\text{max}}) - L(b_{\text{max}})] - 2[L(b_{\text{max}}) - L(b)] + 2[L(b_{\text{max}}) - L(b)]$$

$$\sim \chi^2_{k \text{ of parameters in saturated model}}$$

$$\sim \chi^2_{k \text{ of parameters in model of interest}}$$

\Rightarrow positive constant near zero
 $\&$ model of interest fits data almost as well as the saturated model fits

$$\Rightarrow D \sim \chi^2_{(n-p, r)}$$

\hookrightarrow ncp

\rightarrow Use in tests

\rightarrow D is the TS statistic HT for GLMs

\rightarrow If $y_i \sim \text{Normal} \Rightarrow D \sim \chi^2$ exactly

\rightarrow however it depends on $V(b_{\text{max}}) = \Omega^2$ which is usually unknown

$\Rightarrow D$ cannot be used directly as a test stat

\rightarrow for other dists $\rightarrow D \sim \chi^2_{(n-p, r)}$

\rightarrow however for Binomial & multinom., D can be calculated & used directly as a test stat

\rightarrow Example \rightarrow Let $y_i \sim \text{Bin}(n_i, p_i)$

- Log-likelihood $\rightarrow L(\beta) = \sum_{i=1}^n \left[y_i \ln(\pi_i) + (1-y_i) \ln(1-\pi_i) + w_i \ln(1-T_i) + b_0 + \beta_1 x_i \right]$
- full model \rightarrow MLE \rightarrow for a saturated model, all π_i 's are different $\Rightarrow \beta = (\pi_1, \dots, \pi_n)^T$
- MLE for each β : $\hat{\pi}_i = \hat{y}_i / n_i$
- Likelihood \rightarrow max value of the log-likelihood function is

- Model of interest → ALE → for any other model w/ $p < n$ parameters, let $\hat{\beta}_0$ be the MLE for probabilities π_i w/ $\hat{y}_i = \pi_i x_i$ w/ the fitted values

→ Likelihood → the log likelihood evaluated at those values $\hat{\beta}$

$$L(\hat{\beta}|y) = \left[y_1 \ln\left(\frac{\hat{\pi}_1}{\pi_1}\right) + y_2 \ln\left(\frac{\hat{\pi}_2}{\pi_2}\right) + \dots + y_n \ln\left(\frac{\hat{\pi}_n}{\pi_n}\right) \right]$$

→ Final deviance → $D = -2[L(\text{fitted}) - L(\hat{\beta})]$

$$\downarrow = -2 \sum_{i=1}^n \left[y_i \ln\left(\frac{\hat{\pi}_i}{\pi_i}\right) + (\hat{\pi}_i - \pi_i) \ln\left(\frac{\hat{\pi}_i - \pi_i}{\hat{\pi}_i \pi_i}\right) \right]$$

→ Example → normal model → $E(y_i) = \mu_i = x_i^T \beta$, $y_i \sim N(\mu_i, \sigma^2)$, $i=1, \dots, n$

→ Likelihood → $L(\beta|y) \propto -\frac{1}{2\sigma^2} \sum (y_i - \mu_i)^2 = -\frac{1}{2\sigma^2} \text{Res}(y|\beta)^2$

→ Full model → all μ_i 's can be different $\Rightarrow \dim(\beta) = n \times 1$

→ Maximizing above likelihood gives $\hat{\mu}_i = y_i$

→ Mean value of likelihood \Rightarrow

$$L(\text{fitted}|y) = -\frac{1}{2\sigma^2} N \ln(2\pi\sigma^2)$$

→ Model of interest → let $b = (x^T x)^{-1} x^T y$ be the MLE

→ The corresponding mean value of the likelihood \Rightarrow

$$L(b|y) = -\frac{1}{2\sigma^2} \sum (y_i - b_i)^2 = -\frac{1}{2\sigma^2} \text{Res}(y|b)^2$$

→ Final deviance $\Rightarrow D = -2[L(\text{fitted}) - L(b)]$

$$\downarrow = \frac{1}{\sigma^2} \sum (y_i - b_i)^2$$

$$\downarrow = \frac{1}{\sigma^2} \sum (y_i - \hat{\mu}_i)^2$$

→ If only one parameter $\rightarrow E(y_i) = \mu_i$ for all i , $X = \begin{bmatrix} 1 & \dots & 1 \end{bmatrix}$, $b = \hat{\mu} = \bar{y}$

→ $D = \frac{1}{\sigma^2} \sum (y_i - \bar{y})^2$

→ This is related to the sample variance s^2

$$s^2 = \frac{1}{n-1} \sum (y_i - \bar{y})^2 = \frac{\text{Res}(y|b)^2}{n-1}$$

→ we know $\frac{n-1}{n} s^2 \sim \chi^2_{n-1} \Rightarrow D \sim \chi^2_{n-p}$ exactly

→ More generally $\rightarrow D = \frac{1}{\sigma^2} \sum (y_i - x_i^T b)^2$

$$\downarrow = \frac{1}{\sigma^2} (y - x^T b)^T (y - x^T b)$$

$$\downarrow = y - x(x^T x)^{-1} x^T y$$

$$\downarrow = (I - H)y$$

Wk matrix $= x(x^T x)^{-1} x^T$

→ quadratic form $\rightarrow (y - x^T b)^T (y - x^T b) = \{ (I - H)y \}^T \{ (I - H)y \}$

$$\downarrow = y^T (I - H)y$$

\downarrow rank $H = p$ $\rightarrow H^T = H$

→ $D \sim \chi^2_{n-p}$

→ $b = (X^T X)^{-1} X^T y$

→ $b^T (I - H)x = 0 \Rightarrow D \sim \chi^2_{n-p}$ exactly

→ Tested deviance $\rightarrow \sigma^2 D = \sum (y_i - \hat{\mu}_i)^2$

→ If model fits the data well, then $D \sim \chi^2_{n-p}$.

The expected value for a RV with σ^2 free is $n-p$

→ $E(D) = n-p$

→ This provides an estimate for σ^2 b/c

$$\hat{\sigma}^2 = \frac{\sum (y_i - \hat{\mu}_i)^2}{n-p}$$

→ For normal linear models, some softwares output scaled deviance σ^2 w/ all $\hat{\sigma}^2$ scaled by scale parameter

→ Deviance D is also related to sum of squares of standardized residuals

$$\sum r_i^2 = \frac{1}{\sigma^2} \sum (y_i - \hat{\mu}_i)^2$$

→ provides rough rule of thumb for the overall magnitude of the standardized residuals

→ If the model fits the data well so that $D \sim \chi^2_{n-p}$, you could expect $E(r_i^2) \approx n-p$

→ Example → Poisson model $\rightarrow Y_i \stackrel{\text{ind}}{\sim} \text{Poisson}(\lambda_i)$

→ Likelihood $\rightarrow L(\beta|y) = \prod_i y_i! \ln(\lambda_i) - \lambda_i - \ln(y_i!)$

→ Full model → all λ_i 's are different $\Rightarrow \beta = [\lambda_1, \dots, \lambda_n]^T$

→ MLE $\rightarrow \hat{\lambda}_i = y_i$

→ Mean value $\rightarrow E(\text{fitted}) = \sum \lambda_i \ln(\lambda_i) - \lambda_i - \sum \ln(\lambda_i)$

→ Model of interest → $p < n$ parameters

→ ALE $\rightarrow b$ can be used to calculate estimates $\hat{\lambda}_i$

→ fitted values $\hat{\lambda}_i = b_i^T \ln(x_i) - b_i - \ln(x_i)$

→ Mean value $\rightarrow E(b) = \sum y_i \ln(\lambda_i) - \sum \hat{\lambda}_i - \sum \ln(\lambda_i)$

→ Final deviance $\rightarrow D = -2[L(\text{fitted}) - L(b)]$

$$\downarrow = -2[\sum y_i \ln(\lambda_i/\hat{\lambda}_i) - \sum (\lambda_i - \hat{\lambda}_i)]$$

→ For most models, it can be shown $E(r_i^2) = E(r_i)$

→ $D = 2 \sum_{i=1}^n r_i^2$

$r_i = \frac{\text{observed response value}}{\text{estimated response value}}$

→ This value can be calculated from the data in this case (will be zero if normal, but value D depends on unknown constant σ^2)

→ Can compare this value to $\sigma^2 \sim \chi^2_{n-p}$

→ Small values (relative to df) mean good fit
→ Large values (relative to df) mean bad fit

5.7 → Hypothesis testing

→ Previous approaches about the parameter vector β of length p can be tested using the sampling dist. of the Wald statistic $(\hat{\beta} - \beta)^T J(\hat{\beta} - \beta) \sim \chi^2_p$. Occasionally the score statistic is used $\hat{U}^T \hat{U}^{-1} \hat{U} \sim \chi^2_p$.

→ An alternative approach is to compare the GOF of the nested or hierarchical models (same prob dist + link function, linear component A_0 is a subset of linear component A_1)

→ New approach \rightarrow $H_0: \beta = \beta_0 = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_p \end{bmatrix} \rightarrow \lambda_0$

→ $H_1: \beta = \beta_1 = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_p \end{bmatrix} \rightarrow \lambda_1$, $\beta_1 \neq \beta_0$

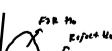
→ We can test H_0 vs H_1 using the difference of deviance stats

1
J

- (although keep in mind practical importance of predictors is not just the statistical significance)

\rightarrow if model M_0 doesn't describe the data well, then $D_0 \stackrel{\text{tend to be}}{\sim} E(D_0)$

\rightarrow if model M_1 , does describe the data well
 $\Rightarrow D_1 \sim \chi^2_{n-p}$ \Rightarrow if M_1 doesn't describe data well,
 $\Rightarrow D_1 \sim \Delta D + \text{larger term expected from } \chi^2_{n-p}$

\rightarrow (usual) testing conclusion \Rightarrow 

\rightarrow If keep in mind model M_0 , may not fit the data particularly well

\rightarrow provided ΔD can be calculated, it is a good method for HT or sampling dist of ΔD (multiple distances) \Rightarrow usually better approximated by χ^2 dist rather than using only a single distance

\rightarrow for normal & some other dists, Deviance isn't totally determined by the set of all subspace parameters that aren't estimated But there can be anomalies \Rightarrow

\rightarrow Example \rightarrow HT Normal Model

$E(y_i) = \mu_i = x_i^T \beta$ $y_i \sim N(\mu_i, \sigma^2)$

$D = \frac{1}{\sigma^2} \sum_i \epsilon_i^2 (\gamma_i - \mu_i)^2$

\rightarrow Here $\mu_1 + \mu_2$ be from M_0 or M_1 respectively

then $D_0 = \frac{1}{\sigma^2} \sum_i \epsilon_i^2 (\gamma_i - \mu_1(x_i))^2$

$D_1 = \frac{1}{\sigma^2} \sum_i \epsilon_i^2 (\gamma_i - \mu_2(x_i))^2$

\rightarrow It's usual to assume that M_1 fits the data well (i.e. it's relevant)

$\Rightarrow D_1 \sim \chi^2_{n-p}$. If M_0 also fits the data well, then $D_0 \sim \chi^2_{n-m}$

$\Rightarrow \Delta D \sim D_0 - D_1$

\rightarrow If M_0 does not fit the data well, (i.e. it's not correct) then $\Delta D \sim \text{noncentral } \chi^2$

To estimate σ^2 , the following ratio is used

$F = \frac{D_0 - D_1}{D_1 / (n-p)}$

$\downarrow = \frac{\left[\sum_i (\epsilon_i^2 (\gamma_i - \mu_1(x_i))^2 - \sum_i (\epsilon_i^2 (\gamma_i - \mu_2(x_i))^2) \right] / (n-p)}{\sum_i (\epsilon_i^2 (\gamma_i - \mu_2(x_i))^2) / (n-p)}$

\rightarrow now F can be directly calculated from the fitted values