

TRANSFORMATIONS AND WEIGHTING TO CORRECT MODEL INADEQUACIES

- Variance Stabilizing Transformations
- Transformations to linearize the model
- Analytical models to select a transformation.
- Generalized & weighted least squares.

The usual approach to deal with inequality of variance is to apply suitable transformation to the response variable or regression variable, i.e., the case when $V(\epsilon_i) \neq \sigma^2$ can be identified from scatter plot of ϵ_i vs. y_i [residual]



Variance Stabilizing Transformation:-

$V(\epsilon) = \sigma^2 \Leftrightarrow$ constant variance assumption.

If the constant variance assumption is violated, the cause is often that the response variable y does not follow a Normal distribution.

Ex. 1.

$$Y \sim \text{Poisson}(\lambda)$$

$$E(Y) = V(Y) = \lambda$$

$Y' = \sqrt{Y}$, then you regress $Y' = \sqrt{Y}$ on x .

$V(\sqrt{Y})$ is independent of mean λ .

Ex. 2.

y is a proportion $0 \leq y \leq 1$

$$y' = \sin^{-1}(\sqrt{y})$$

Constant variance assumption is violated (double-bow pattern in residual plot)

y has mean μ and variance σ^2

situation:- $g(\mu) = \sigma^2$ [Variance depends on mean (not constant)]

Take the following transformation: $U = f(Y) = f(\mu) + \frac{f'(\mu)}{\mu}(Y-\mu)$ [From Taylor's series (next pg.)]

$$V(U) = V(f(Y)) = [f'(\mu)]^2 V(Y) = [f'(\mu)]^2 g(\mu)$$

If we choose the function f such that [since $g(\mu) = \sigma^2$]

$$[f'(\mu)]^2 = \frac{1}{g(\mu)} \Rightarrow f'(\mu) = [g(\mu)]^{-\frac{1}{2}}$$

$$\text{Then } V(U) = V(f(Y)) = 1.$$

$\therefore f(Y)$ is a transformation of Y for which variance becomes constant and this solves the problem on deciding about suitable variance stabilizing transformation for a given problem.

2

Example:

$$\sigma^2 \leq k\mu^q$$

$$g(\mu) = k\mu^q$$

We want $f'(y) = [g(y)]^{-1/2}$

$$\propto \mu^{-q/2}$$

Thus $f(y) \propto y^{1-q/2}$ if $q \neq 2$
 $\log y$ if $q=2$

Commonly used transformations:Relationship of σ^2 and $E(Y)$

$$\sigma^2 \propto \text{constant}$$

$$\sigma^2 \propto E(Y) \quad [q=1]$$

$$\sigma^2 \propto [E(Y)]^2 \quad [q=2]$$

$$\sigma^2 \propto [E(Y)]^3 \quad [q=3]$$

$$\sigma^2 \propto [E(Y)]^4 \quad [q=4]$$

$$\sigma^2 \propto E(Y)(1-E(Y))$$

[when Y is a proportion
between 0 and 1]

Transformation to make constant variance

No transformation

$$f(y) = \sqrt{y} \quad [Y \sim \text{Poisson}]$$

$$f(y) = \log y \quad [Y \sim \text{exponential}]$$

$$f(y) = y^{-1/2}$$

$$f(y) = 1/y$$

$$f(y) = \sin^{-1}(\sqrt{y})$$

Recall Taylor Series expansion:

Taylor series of a real/complex valued function $f(x)$ that is infinitely differentiable function in a neighborhood of a real/complex number 'a' is

$$f(x) = f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f'''(a)}{3!}(x-a)^3 + \dots$$

[3]

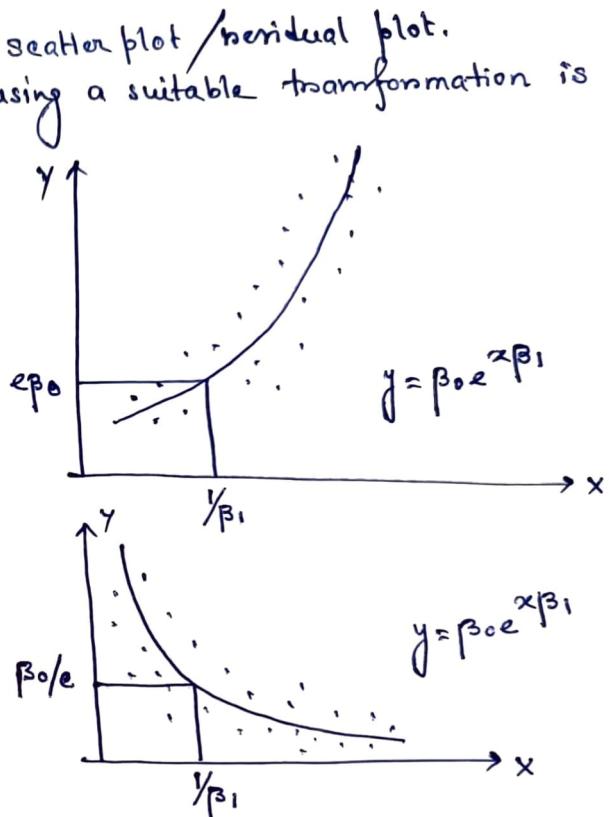
If the relationship between response variable and the regressors variables is NOT linear, we apply transformation to linearize the model:

Transformations to linearize the model:

- Nonlinearity may be detected via scatter plot / residual plot.
- A function can be linearized by using a suitable transformation is called linearizable function.

Example 1: If the scatter plot of y on x suggest an exponential relationship between them, then the appropriate model will be

$$y = \beta_0 e^{x\beta_1}$$



- This model is linear because this is equivalent to the model

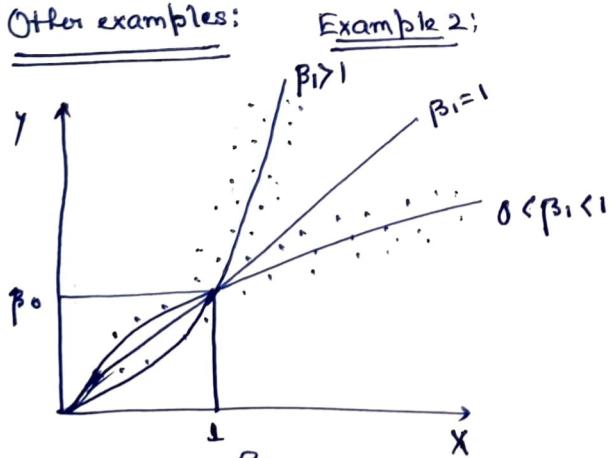
$$\log y = \log \beta_0 + \beta_1 x$$

- Transformation: $y' = \log y$

$$y' = \beta_0' + \beta_1 x$$

- Fit a linear model using $(\log y, x)$ transformed data.

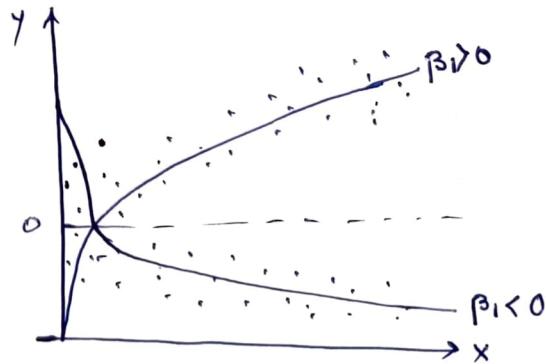
Other examples:



Example 2:

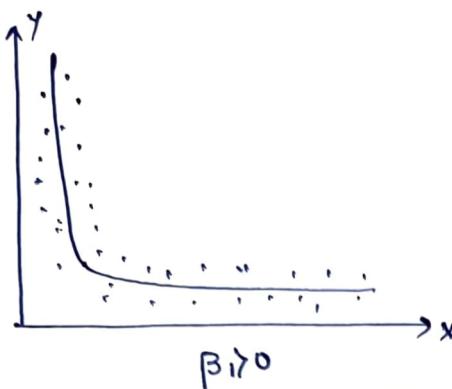
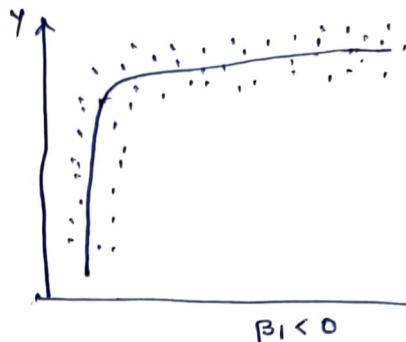
- Model: $y = \beta_0 x^{\beta_1}$
- $\log y = \log \beta_0 + \beta_1 \log x$
- Transformation: $y' = \log y$ and $x' = \log x$
- $y' = \beta_0' + \beta_1 x'$
- Fit a linear model using $(\log y, \log x)$ transformed data.

Example 3:



- $y = \beta_0 + \beta_1 \log x \quad (\beta_1 > 0)$
- Transformation: $x' = \log x$
 $\Rightarrow y = \beta_0 + \beta_1 x'$
- Fit a linear model with $(\log x, y)$.

A

Example 4:Example 5:

$$y = \frac{x}{\beta_0 x - \beta_1}$$

$$\frac{1}{y} = \frac{\beta_0 x - \beta_1}{x} = \beta_0 - \frac{\beta_1}{x}$$

Transformation:

$$y' = 1/y, x' = 1/x$$

$$\Rightarrow y' = \beta_0 - \beta_1 x'$$

Thus, in both the examples, we can fit a linear model using (x', y') transformed data.

Box-Cox METHOD: A useful class of transformation is power transformation y^λ , where λ is a parameter to be determined. Major disadvantage of this method is: as $\lambda \rightarrow 0$, $y^\lambda \rightarrow 1$, all the response values are equal to 1.

To solve this difficulty, Box-Cox method is used:

$\lambda \uparrow, w \uparrow$ thus it is impractical to compare neg. models

$\lambda \approx \text{large} \Rightarrow w \approx \text{large}$

To solve this

$$w = \begin{cases} \frac{y^\lambda - 1}{\lambda} & \text{for } \lambda \neq 0 \\ \log y & \text{for } \lambda = 0 \end{cases}$$

[Assuming y values to be positive]

$$v = \begin{cases} \frac{(y + \lambda_1)^{\lambda_2} - 1}{\lambda_2} & \text{if } \lambda_2 \neq 0 \\ \log(y + \lambda_1) & \text{if } \lambda_2 = 0 \end{cases}$$

[Assuming y values to be negative]

Geometric mean: $\tilde{Y} = \left(\prod_{i=1}^n y_i \right)^{1/n}$ is used as normalization factor.

Given data: (x_i, y_i)

$$w' = \begin{cases} \frac{y^\lambda - 1}{\lambda \tilde{Y}^{\lambda-1}} & \text{for } \lambda \neq 0 \\ \tilde{Y} \log y & \text{for } \lambda = 0 \end{cases}$$

$(y, \dots, y_n) \rightarrow (w'_1, w'_2, \dots, w'_n)$
and use $w' = X\beta + \epsilon$ by OLS
for any specified value of λ .

Maximum Likelihood Method for estimating λ :

- Choose a value of λ from $[-2, 2]$ at first and extend the range later, if necessary.
- For each chosen λ value, estimate λ and compute $SS_{Reg}(\lambda)$ for the regression model $w' = X\beta + \epsilon$.
- The MLE of λ corresponds to the value of λ for which $SS_{Reg}(\lambda)$ is minimum.

WEIGHTED LEAST SQUARES METHOD:

Linear regression model with "non-constant variance" can be fitted by the method of weighted least squares (WLS).

Simple linear regression model: $y = \beta_0 + \beta_1 x + \epsilon$ $[y_i \sim N(\beta_0 + \beta_1 x_i, \sigma_i^2)]$

The weighted least square function is $S = \sum w_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 = \sum w_i e_i^2$
and $w_i \propto \frac{1}{\sigma_i^2}$

Normal equations are:

$$\frac{\partial S}{\partial \beta_0} = 0 \Rightarrow \sum w_i y_i = \hat{\beta}_0 \sum w_i + \hat{\beta}_1 \sum w_i x_i$$

$$\frac{\partial S}{\partial \beta_1} = 0 \Rightarrow \sum w_i y_i x_i = \hat{\beta}_0 \sum w_i x_i + \hat{\beta}_1 \sum w_i x_i^2$$

Question: Why $w_i \propto \frac{1}{\sigma_i^2}$ and how to know σ_i^2 ? [NEXT PG]

GENERALISED LEAST SQUARES METHOD:

- Weighted least square for multiple regression model.
- Consider the same model $y = X\beta + \epsilon$ with $E(\epsilon) = 0$ and $V(\epsilon) = V\sigma^2$ and $V(\epsilon) \neq I\sigma^2$ (constant variance), where V is a positive definite matrix.
- This happens when:
 - Observations y have unequal variances and/or
 - observations are correlated.
- In either case, the condition of Gauss-Markov theorem are violated. So, $\hat{\beta} = (X'X)^{-1}X'y$ is NOT BLUE.
- However, it is possible to find BLUE of β for arbitrary positive definite V by suitable linear transformation on the model.

$$(X, Y) \rightarrow (GX, GY)$$

$$GY = GX\beta + G\epsilon$$

$$V(G\epsilon) = \sigma^2 GVG'$$
 since $V(\epsilon) = V\sigma^2$.

$$V(G\epsilon) \text{ is constant if } GVG' = I. \text{ Then } V(G\epsilon) = \sigma^2 I.$$

Therefore if we choose G such that $GVG' = I$ then the transformed data satisfy Gauss-Markov conditions and the BLUE of β is obtained by OLS estimation of transformed data.

$$GVG' = I \Leftrightarrow V^{-1} = G'G$$

$$\Leftrightarrow V = G^{-1}G'^{-1}$$

$$(X, Y) \rightarrow (GX, GY)$$

$$\left| \begin{array}{l} \hat{\beta} = (X'G'GX)^{-1}X'G'GY \\ = (X'V^{-1}X)^{-1}X'V^{-1}Y \\ V(\hat{\beta}) = \sigma^2 (X'G'GX)^{-1} = \sigma^2 (X'V^{-1}X)^{-1}. \end{array} \right.$$

- It is always possible to find a symmetric G_1 by using the orthogonal decomposition of positive definite V

$$V = U' \Lambda^{-1} U,$$

where Λ is the diagonal matrix of eigenvalues and U is the matrix of eigenvectors.

$$V^{-1} = U^{-1} \Lambda U'^{-1}$$

$$\Leftrightarrow V^{-1} = G_1' G_1 ; \text{ where } G_1 = \Lambda^{\frac{1}{2}} U'^{-1} \quad [\text{choice of } G_1]$$

- Thus, if the data does not satisfy the assumption of constant variance and if the errors/observations are correlated then $V(\epsilon) = V\sigma^2$, then we can find a transformation of the data $(G_1 X, G_1 Y)$ which satisfy all the assumptions and we can apply OLS technique on the transformed data.
- WLS technique as a particular case of GLS:

- Observations are uncorrelated and have unequal variance.

$$V(\epsilon) = \begin{pmatrix} \sigma_1^2 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & & \vdots \\ 0 & 0 & \dots & \sigma_n^2 & \end{pmatrix} = V$$

$$V^{-1} = G_1' G_1 \Rightarrow G_1 = V^{-1/2} = \begin{pmatrix} 1/\sigma_1 & 0 & \dots & 0 \\ 0 & 1/\sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & 1/\sigma_n \end{pmatrix}$$

$$= \text{Diag}\left(\frac{1}{\sigma_1}, \frac{1}{\sigma_2}, \dots, \frac{1}{\sigma_n}\right); W = V^{-1}$$

$$w_i \propto \frac{1}{\sigma_i^2}$$

$$\hat{\beta} = (X' V^{-1} X)^{-1} X' V^{-1} Y$$

$$= (X' W X)^{-1} (X' W Y)$$

$$\text{Var}(\hat{\beta}) = (X' V^{-1} X)^{-1} = (X' W X)^{-1}.$$

- GLS technique is used to deal with non-constant variance in the response variable or observations are correlated.

7

EXAMPLE:

Suppose we have n observations of variables X_1, X_2, \dots, X_k, Y , where X_i 's are predictor variables and Y is a response variable. Suppose we are told that observations Y_i are uncorrelated but the last observation has variance $4\sigma^2$ rather than σ^2 . Find the best linear unbiased estimator (BLUE) of β using weighted least square?

Solution:

$$V(Y_i) = V(\epsilon_i) = \sigma^2, i=1(1)n-1, V(Y_n) = V(\epsilon_n) = 4\sigma^2$$

$$V(\epsilon) = \begin{pmatrix} \sigma^2 & 0 & 0 & \cdots & 0 \\ 0 & \sigma^2 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 4\sigma^2 \end{pmatrix} = V(\sigma^2) \text{ where } V = \text{diag}(1, 1, \dots, 1, 4).$$

$$Y = X\beta + \epsilon \quad \text{where } \text{Var}(Y_i) = V\sigma^2$$

$$G_1 Y = G_1 X \beta + G_1 \epsilon$$

$$V(G_1 \epsilon) = \sigma^2 G_1 V G_1' = \sigma^2 I$$

$$G_1 V G_1' = I$$

$$\Rightarrow V^{-1} = G_1' G_1.$$

$$\hat{\beta} = (X' G_1' G_1 X)^{-1} X' G_1' G_1 Y \text{ is BLUE}$$

$$= (X' V^{-1} X)^{-1} X' V^{-1} Y$$

$$= (X' \text{diag}(1, \dots, 1, 1/4) X)^{-1} (X' \text{diag}(1, 1, \dots, 1/4) Y).$$

Generalized Linear Model (GLM)

- GLM analysis comes into account when errors distribution is not normal (response is not normal) and/or (but must be a member of exponential family) when a vector of non-linear functions of the responses $\eta(Y) = (\eta(Y_1), \eta(Y_2), \dots, \eta(Y_n))'$ are not Y itself, has expectation $X\beta$.

In standard MLR :- $Y = X\beta + \epsilon$

$$E(Y) = X\beta, E(\epsilon) = 0$$

\therefore Here we assume

$$E(\eta(Y)) = X\beta$$

- In GLM, the response variable distribution must be a member of the exponential family.
- The exponential family of distribution:

A random variable u that belongs to exponential family with single parameter θ has a pdf

$$f(u, \theta) = s(u) t(\theta) e^{a(u)b(\theta)},$$

where s, t, a, b are all known functions.

Rewrite: $f(u, \theta) = \exp \{a(u)b(\theta) + d(u) + c(\theta)\},$

where $d(u) = \ln(s(u))$, $c(\theta) = \ln(t(\theta))$ and

when $a(u) = u$, the distn. is said to be in canonical form.

Note: $b(\theta)$ is called natural parameter, parameters other than the parameter of interest (θ) are called nuisance parameters.

- Some members of exponential family:

1 Normal distribution: - $N(\mu, \sigma^2)$ and parameter of interest: μ .

$$f(u, \mu) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(u-\mu)^2}{\sigma^2}}, -\infty < u < \infty$$

$$= \exp \left\{ u \cdot \frac{\mu}{\sigma^2} + \left[-\frac{\mu^2}{2\sigma^2} - \frac{1}{2} \ln 2\pi\sigma^2 \right] - \frac{u^2}{2\sigma^2} \right\}$$

where, $a(u) = u$, $b(\theta) = \frac{u}{\sigma^2}$, $c(\theta) = -\frac{\mu^2}{2\sigma^2} - \frac{1}{2} \ln 2\pi\sigma^2$, $d(u) = \frac{u^2}{2\sigma^2}$.

2 Binomial distribution: - $\text{Bin}(n, p)$ and p is the parameter of interest and n is the nuisance parameter

$$f(u, p) = \binom{n}{u} p^u (1-p)^{n-u}, \quad u = 0, 1, \dots, n$$

$$= \binom{n}{u} \left(\frac{p}{1-p} \right)^u (1-p)^n$$

$$= \exp \left\{ u \ln \left(\frac{p}{1-p} \right) + n \ln(1-p) + \ln \binom{n}{u} \right\}$$

where, $a(u) = u$, $b(\theta) = \ln \left(\frac{p}{1-p} \right)$, $c(\theta) = n \ln(1-p)$, $d(u) = \ln \binom{n}{u}$.

Similarly, we can show Poisson distribution, Gamma distribution, exponential distribution, Negative Binomial distn., etc. belong to exponential family of distributions.

Expected value and variance of $a(u)$:

$$E(a(u)) = -\frac{c'(\theta)}{b'(\theta)}, V(a(u)) = \frac{b''(\theta) c'(\theta) - c''(\theta) b'(\theta)}{[b'(\theta)]^3} \quad (*)$$

- Fitting GLM :- Suppose we have a set of independent obsn.s. (Y_i, \tilde{x}_i') , $i=1(1)n$, $\tilde{x}_i' = (x_{i1}, x_{i2}, \dots, x_{ip})$ from some exponential type distribution of canonical form [i.e., $a(Y) = Y$]. The joint pdf is

$$\{f(Y_1, Y_2, \dots, Y_n, \theta, \phi) = \exp \left\{ \sum_{i=1}^n Y_i b(\theta) + \sum_{i=1}^n c(\theta_i) + \sum_{i=1}^n d(Y_i) \right\},$$

where ϕ is a vector of a nuisance parameters that occur with in $b(\cdot)$, $c(\cdot)$, & $d(\cdot)$.

Hence $\theta = (\theta_1, \theta_2, \dots, \theta_p)$: vector of parameters of interest.

- The variation in response variable (Y_i) can be explained in terms of \tilde{x}_i values, i.e., $\tilde{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})'$
- Consider the set of parameters $\beta = (\beta_1, \beta_2, \dots, \beta_p)'$ (vector of regression coefficients)
We find some suitable link function $g(\cdot)$ such that

$$g(\mu_i) = \tilde{x}_i' \beta \quad | \quad y_i = \tilde{x}_i' \beta + \epsilon \quad \begin{matrix} \text{(standard MLR)} \\ Y_i \sim \text{Normal} \\ g(\theta) = \theta \end{matrix}$$

- A link function that is often regarded as sensible one is natural parameter.

Example:- Binomial Distribution

Suppose we have data (Y_i, \tilde{x}_i') from a binomial distn. $\text{Bin}(n_i, p_i)$.
The single observation Y_i is of the form $\frac{r_i}{n_i}$, where r_i is the no. of success in n_i trials, each having prob. p_i of success (parameter of interest).
and $\tilde{x}_i' = (x_{i1}, x_{i2}, \dots, x_{ip})$ is a set of observations of p regressors associated with Y_i . Binomial distn. is a member of exp. family: (n_i is the nuisance parameter).

[*] Example: For Binomial distribution: $a(\theta) = u, b(\theta) = \ln(\frac{p}{1-p}), c(\theta) = n \ln(1-p)$

$$E(a(u)) = E(u) = -\frac{c'(\theta)}{b'(\theta)} = \frac{n}{1-p} \times p(1-p) = np.$$

$$Var(a(u)) = V(u) = \frac{b''(\theta) c'(\theta) - c''(\theta) b'(\theta)}{[b'(\theta)]^3} = np(1-p).$$

$$\begin{aligned}
 \text{Joint pdf} &= f(Y_1, Y_2, \dots, Y_n) \\
 &= \prod_{i=1}^n \left(\frac{n_i}{Y_i} \right)^{Y_i} p_i^{Y_i} (1-p_i)^{n_i - Y_i} \\
 &= \prod_{i=1}^n \exp \left\{ Y_i \ln \left(\frac{p_i}{1-p_i} \right) + n_i \ln (1-p_i) + \ln \left(\frac{n_i}{Y_i} \right) \right\} \\
 &= \exp \left\{ \sum_{i=1}^n Y_i \ln \left(\frac{p_i}{1-p_i} \right) + \sum_{i=1}^n n_i \ln (1-p_i) + \sum_{i=1}^n \ln \left(\frac{n_i}{Y_i} \right) \right\}
 \end{aligned}$$

$$\text{Natural parameter} = \ln \left(\frac{p_i}{1-p_i} \right).$$

We would hope that the variation in the $Y_i / E(Y_i) = p_i$ could be explained in terms of \tilde{x}_i' values, i.e., we would hope that we could find a suitable link function $g(\cdot)$ such that $g(p_i) = \tilde{x}_i' \beta$.

We fit the model $\ln \left(\frac{p_i}{1-p_i} \right) = \tilde{x}_i' \beta = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$

$$E(Y_i) = p_i = \frac{\exp(\tilde{x}_i' \beta)}{1 + \exp(\tilde{x}_i' \beta)}$$

In stead of fitting $y_i = \tilde{x}_i' \beta + \epsilon$ we fit here

$$y_i = p_i + \epsilon,$$

where $\tilde{x}_i' \beta = \beta_1 + \beta_2 x_{i2}$ (*) is called the logistic function.

Estimation via ML function: To estimate β , we use maximum likelihood method:

$$\text{Likelihood func: } L = \exp \left\{ \sum_{i=1}^n y_i \ln \left(\frac{p_i}{1-p_i} \right) + \sum_{i=1}^n n_i \ln (1-p_i) + \sum_{i=1}^n \ln \left(\frac{n_i}{y_i} \right) \right\}$$

$$\ln L = \sum_{i=1}^n y_i \ln \left(\frac{p_i}{1-p_i} \right) + \sum_{i=1}^n n_i \ln (1-p_i) + \sum_{i=1}^n \ln \left(\frac{n_i}{y_i} \right)$$

$$= \sum y_i \tilde{x}_i' \beta - \sum n_i \ln (1 + \exp(\tilde{x}_i' \beta)) + \sum \ln \left(\frac{n_i}{y_i} \right)$$

Maximize $\ln L$ w.r.t. β , use numerical search/iteratively reweighted least square (IRLS) could be used to

compute MLE of β .
Choice of Link function:-

Normal : $g(\mu) = \mu$ (Identity link)

Binomial : $g(p) = \ln \left(\frac{p}{1-p} \right)$ (logistic link)

Poisson : $g(\mu) = \ln \mu$ (log link)

Gamma/Exponential : $g(\mu) = \frac{1}{\mu}$ (reciprocal link)

Example (Pneumoconiosis Data)
Lung disease result from breathing & dust in coalmines

Number of years of exposure (x_i)	Number of severe cases	Total number of miners	Proportion of severe cases (y_i)
5.8	0	98	0
15.0	1	54	0.0185
21.5	3	43	0.0698
27.5	8	48	0.1667
33.5	9	51	0.1765
39.5	8	38	0.2105
46.0	10	28	0.3571
51.5	5	11	0.4545

Check the variation in y_i can be explained in terms of the number of years of exposure?

Sol. Probability distr. for the number of severe cases is binomial.
We will fit a logistic regression model to the data.

$$p_i = E(y_i) = \frac{\exp(\tilde{x}_i' \beta)}{1 + \exp(\tilde{x}_i' \beta)}$$

Here $\tilde{x}_i' \beta = \beta_1 + \beta_2 x$, where x : # years of exposure.

$$\hat{y}_i = \frac{\exp(4.79 - 0.0935x)}{1 + \exp(4.79 - 0.0935x)}.$$

Poisson Regression Model:

Data (y_i, \tilde{x}_i') from $P(\mu_i)$, $E(y_i) = \mu_i$.

$f(y, \mu) = \exp\{y \ln \mu - \mu - \ln y!\}$ where $\ln \mu$ is the natural parameter,
i.e., $g(\mu) = \ln \mu$ is the link function.

The variation in y_i could be explained in terms of the \tilde{x}_i' values. We fit the model: $g(\mu_i) = \tilde{x}_i' \beta$

$$\ln \mu_i = \tilde{x}_i' \beta = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}.$$

$\mu_i = e^{\tilde{x}_i' \beta}$; final model to fit.

$E(y_i) = e^{\tilde{x}_i' \beta}$. Same as writing $y_i = e^{\tilde{x}_i' \beta} + \epsilon$ [Final model]

In SLR: $y = \beta_0 + \beta_1 x + \epsilon$; Assumption: $\epsilon \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$

ϵ : random variable

x_i : controlled variable/deterministic variable.

Two variations of situation could be:

(a) Response and the regression variables are jointly distributed RVs.

(b) There are measurement errors in regressors.

Measurement errors & Calibration Problem

- Case where Response & Regressors are jointly distributed RVs.
- Measurement Errors in Regressors
- The Calibration problem (inverse problem)

① $X \text{ & } Y$ are jointly normally distributed:

$$f(x, y) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp \left\{ -\frac{1}{2(1-\rho^2)} \left[\left(\frac{y-\mu_1}{\sigma_1} \right)^2 + \left(\frac{x-\mu_2}{\sigma_2} \right)^2 - 2\rho \left(\frac{y-\mu_1}{\sigma_1} \right) \left(\frac{x-\mu_2}{\sigma_2} \right) \right] \right\}$$

$$E(Y) = \mu_1, V(Y) = \sigma_1^2, E(X) = \mu_2, V(X) = \sigma_2^2$$

and $\rho = \frac{E[(Y-\mu_1)(X-\mu_2)]}{\sigma_1\sigma_2}$ is correlation coefficient between Y and X .

∴ The conditional distn. of Y given X is

$$Y|X \sim N\left(\mu_1 + \rho \frac{\sigma_1}{\sigma_2}(X - \mu_2), (1 - \rho^2)\sigma_1^2\right).$$

$$E(Y|X) = \mu_1 + \rho \frac{\sigma_1}{\sigma_2}(X - \mu_2)$$

$$E(Y|x) = \beta_0 + \beta_1 x ; \quad \text{where } \beta_0 = \mu_1 - \rho \mu_2 \frac{\sigma_1}{\sigma_2}, \quad \beta_1 = \rho \frac{\sigma_1}{\sigma_2}$$

Model when X is a deterministic variable $\begin{cases} E(Y) = \beta_0 + \beta_1 x \\ Y = \beta_0 + \beta_1 x + \epsilon; \quad \epsilon \sim N(0, \sigma^2) \end{cases}$ when both X and Y are RVs.

$$Y_i | x_i \stackrel{\text{ind.}}{\sim} N(\beta_0 + \beta_1 x_i, (1 - \rho^2)\sigma_1^2) \quad [\text{Model: } E(Y|x) = \beta_0 + \beta_1 x]$$

$$\text{MLE: } L = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_1^2(1-\rho^2)}} e^{-\frac{1}{2\sigma_1^2(1-\rho^2)}(y_i - \beta_0 - \beta_1 x_i)^2}$$

$$= \left(\frac{1}{\sqrt{2\pi\sigma_1^2(1-\rho^2)}} \right)^n e^{-\frac{1}{2\sigma_1^2(1-\rho^2)} \sum (y_i - \beta_0 - \beta_1 x_i)^2} \quad [\text{Max } L \equiv \min \sum (y_i - \beta_0 - \beta_1 x_i)^2]$$

We find β_0 & $\beta_1 \Rightarrow \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$ is minimum.

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \text{ and } \hat{\beta}_1 = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}} = \sqrt{\frac{S_{ST}}{S_{xx}}} \cdot n$$

We know (from OLS) identical to those given by LSE in case where x is a controlled variable.

$$\rho = \text{Corr}(X, Y); \text{ estimate of } \rho \text{ is } r = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}} = \text{sample correlation coefficient.}$$

To test $H_0: \rho = 0$ vs $H_1: \rho \neq 0$: Test statistic is $t_0 = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$ ~ t_{n-2} under H_0 .

We Reject H_0 if $|t_0| > t_{\alpha/2, n-2}$. When $\rho = 0$, there is no relationship between X and Y .

Measurement errors in Regressions:-

We wish to fit the simple linear reg. model, but the regressors is measured with errors.

$$x_i = \tilde{x}_i + a_i \quad ; \quad E(x_i) = \tilde{x}_i \quad ; \quad E(a_i e_i) = 0$$

↓ ↓ ↓
Observed value true value measurement errors

Assume: $E(a_i) = 0$ and $V(a_i) = \sigma_a^2$. Also, a_i and e_i are independent.

The response variable is subject to the usual errors ϵ_i $(i=1(1)n)$.

The reg. model is $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ [Given data: (y_i, x_i)]

$$\begin{aligned} &= \beta_0 + \beta_1 (\tilde{x}_i - a_i) + \epsilon_i \\ &= \beta_0 + \beta_1 x_i + (\epsilon_i - \beta_1 a_i) \end{aligned}$$

Since x_i and y_i both are RVs, then we need to compute, $Cov(x_i, y_i)$

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad ; \quad \epsilon_i = \epsilon_i - \beta_1 a_i$$

$$\begin{aligned} Cov(x_i, y_i) &= E[(x_i - E(x_i))(\epsilon_i)] \\ &= E[(\tilde{x}_i - a_i)(\epsilon_i - \beta_1 a_i)] \\ &= E[a_i(\epsilon_i - \beta_1 a_i)] \\ &= -\beta_1 E(a_i^2) \\ &= -\beta_1 \sigma_a^2 \end{aligned}$$

If we apply standard LSM to the data, the estimates of the model parameter are no longer unbiased.

Usual case: $\hat{\beta}_1 = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$

$$\hat{\beta}_1 = \frac{\beta_1}{1 + \theta} \quad \left[\text{Hence in the presence of measurement error} \right]$$

$$\text{where } \theta = \frac{\sigma_a^2}{\sigma_x^2}$$

$$\sigma_x^2 = \frac{\sum (x_i - \bar{x})^2}{n}$$

$\hat{\beta}_1$ is a biased estimator of β_1 , unless $\sigma_a^2 = 0$, that is there is no measurement error in regressions.

- If σ_a^2 is small relative to σ_x^2 , the bias will be small [$\theta \approx 0$].
- If variability in the measurement errors is small relative to the variability of the x 's; then measurement error can be ignored & OLS method can be applied directly.

Example (Measurement Error Problem)

X : current flow in an electric circuit
 Current flow is measured with an ammeter which is not completely accurate. Hence, a measurement error is experienced.

$$X = x + \epsilon$$

\downarrow

true current flow

The Calibration Problem: →

- Given an observed value of y , say y_0 , determine the x value corresponding to it. (Inverse problem)
- Example: We know the temperature reading given by thermocouple is a linear function of the actual temperature.

$$\text{Observed temp} = \beta_0 + \beta_1 (\text{actual temp}) + \epsilon$$
- Solution:- Suppose we have $(y_i, x_i), i=1(1)n$.
 First, fit the model $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$.
 Let y_0 be the observed value of y . A natural point estimation of the corresponding value of x is $\hat{x}_0 = \frac{y_0 - \hat{\beta}_0}{\hat{\beta}_1}$, assuming that $\hat{\beta}_1 \neq 0$.