

## Topic 1

### Review of Linear Regression/Normal Error Regression Model:

- Build and fit a model which describes the relationship between a dependent response variable and set of explanatory or predictor variables
- $Y_i = \alpha_0 + \alpha_1 * X_i + \varepsilon_i$
- $Y_i$  is the  $i^{\text{th}}$  observation of dependent variable  $Y$ ,  $X_i$  is the  $i^{\text{th}}$  observation of independent variable  $X$ ,  $\alpha_0$  and  $\alpha_1$  are unknown parameters and  $\varepsilon_i$  is a random error term which is  $N(0, \sigma^2)$
- $\sigma^2$  is constant, not dependent on  $X_i$  and  $Cov(\varepsilon_i, \varepsilon_j) = 0$ , uncorrelated error
- Response  $Y_i$  is decomposed into systematic (non-random term):  $\alpha_0 + \alpha_1 X_i$  and non-systematic (random term):  $\varepsilon_i$
- $E[Y_i] = \alpha_0 + \alpha_1 X_i$  and  $Var[Y_i] = \sigma^2$
- Normality Assumption:  $\varepsilon_i$  are normally distributed and so the response  $Y_i$  is normally distributed

### Estimation of $\alpha_0$ and $\alpha_1$ :

- Can be done under OLS or MLE
- $\hat{\alpha}_0 = \bar{Y} - \hat{\alpha}_1 * \bar{X}$
- $\hat{\alpha}_1 = \frac{\sum(X_i - \bar{X})(Y_i - \bar{Y})}{\sum(X_i - \bar{X})^2}$
- Under normality assumptions for  $\varepsilon_i$ , OLS and MLE methods produce the same estimates

### Fitted Regression Line

- Obtain  $\hat{\alpha}_0$  and  $\hat{\alpha}_1$
- The fitted line equation is  $\hat{Y}_i = \hat{\alpha}_0 + \hat{\alpha}_1 * X_i$
- $\hat{Y}_i$  is the estimate of the mean response  $E[Y_i]$
- Regression line is a straight line in the X-Y plane
- Residuals:  $\hat{\varepsilon}_i = Y_i - \hat{Y}_i$
- $\hat{\varepsilon}_i$  is an estimate of  $\varepsilon_i$
- MLE of  $\sigma^2$  is the mean of the squared residuals
- $\hat{\sigma}^2 = \frac{1}{n} * \sum \hat{\varepsilon}_i^2$ , this is a biased estimate
- Unbiased estimate of  $\sigma^2$  is  $\frac{1}{n-2} * \sum \hat{\varepsilon}_i^2$ , subtract 2 degrees of freedom due to 2 estimations of alpha's

### Testing Assumptions:

- Use residuals to test if the regression is linear, if random error terms have constant variance, if random error terms are independent and if error terms are Normal Dist
- If some assumptions are not valid, regression model doesn't fit the data well

### Intro to GLMs:

- Normal error linear regression model may not be suitable as assumptions underlying the model are not appropriate
- A generalised linear model (GLM) is more flexible since it incorporates the non-normality of the response and nonlinear relationship between mean response  $E[Y_i]$  and predictor  $X_i$ 's

- Detect normal linear model is not suitable through common sense, histogram of the data, a normal probability plot (Q-Q plot) or a normality test (Anderson darling)
- Q-Q plot = Need straightish line, density = symmetrical
- Detect non-constancy of variance through data plot or residual plot
- Implication of non-constancy – mean and variance are totally unrelated to each in the Normal distribution but for many non-normal distributions, the variance is linked to the mean. Hence the variance is not constant as it changes with mean, and the normality assumption for the response  $Y_i$  or  $\varepsilon_i$  is not valid and therefore normal error regression model is not likely to provide a good fit.
- Under a GLM, response  $Y_i$  can be any distribution from the exponential family (Normal, Exp, Poisson, Gamma, Binomial, Inverse Gaussian, Neg Bin)
- Exponential family has good mathematic properties to estimate the unknown parameters easily using the maximum likelihood estimation

### Link Functions and Examples:

- In normal error regression model, mean of response is a linear function of the explanatory variables, if  $\mu_i = E(Y_i)$ , then
 
$$\mu_i = \beta_0 + \beta_1 * X_{1i} + \dots + \beta_n * X_{ni}$$
- In a GLM, the relationship between the mean response and the explanatory variables may not be linear, can be exp fn, reciprocal etc
- Linear Predictor  $\eta_i$
- $\eta_i = \beta_0 + \beta_1 * X_{1i} + \dots + \beta_n * X_{ni}$
- There exists a relationship between the mean response  $\mu_i$  and the predictors which can be expressed as  $\mu_i = \eta_i$  or  $\mu_i = \exp(\eta_i)$  etc
- These transformations are usually rearranged to be a function of  $\mu_i$  and are called link functions
- $\eta_i = \mu_i$  is the identity link function
- $\eta_i = \ln(\mu_i)$  is the log link function
- $\eta_i = \frac{1}{\mu_i}$  is the reciprocal link function
- $\eta_i = \ln\left(\frac{\mu_i}{1-\mu_i}\right)$  is the logit link function
- Generally, if  $\eta_i = g(\mu_i)$ ,  $g(\mu)$  is called the link function

### Exponential Family:

- For GLM, assume the distribution of the random error term is in the exponential family (Exp, Bernoulli, Binomial, Normal, Poisson, Gama, NB)
- A distribution belongs to the exp family if it can be written in the form

$$f(y; \theta, \phi) = \exp\left(\frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi)\right)$$

- a,b,c are function
- $\theta$  is the natural parameter and  $\phi$  is the scale parameter
- Normal:  $\theta = \mu, \phi = \sigma^2, b(\theta) = \frac{\theta^2}{2}, a(\phi) = \phi, c(y, \phi) = -\frac{y^2}{2\phi} - \frac{1}{2}\ln(2\phi\pi)$
- Poisson:  $\theta = \ln(\lambda), \phi = 1, b(\theta) = \exp(\theta), a(\phi) = 1, c(y, \phi) = -\ln y!$
- Method: Put everything into exp using  $e^{\ln \cdot}$  and rearrange to get form
- When find  $\theta$ , sub back into  $f()$  to get all others