

LOGISTIC REGRESSION ASSUMPTIONS

- When the assumptions of logistic regression analysis are not met problem such as biased coefficient may lead to invalid statistic inference.

ASSUMPTION 1 → Must be Output Categorical Column.

ASSUMPTION 2 → Linearity of independent variable and log odds.

ASSUMPTION 3 → No strong influential outliers (Cook distance, Box plot)

ASSUMPTION 4 → Absence of Multicollinearity (Correlation & VIF)

ASSUMPTION 5 → Independence of Observations (Autocorrelation, Residual plot)

ASSUMPTION 6 → Sufficiently large sample size

ASSUMPTION 1 → Output is Categorical Column.

→ We can directly check datatype and distribution of the target column.

ASSUMPTION 2 → Linearity of independent variable and log odds.

→ Log of Odds solve the problem of symmetry.

- Suppose my team is bad

$$\text{Odd} = \frac{1}{4} = 0.25, \text{ Worse} = \frac{1}{8} = 0.125$$

$$\text{More worse} = \frac{1}{16} = 0.06, \text{ Full worse} = \frac{0}{\text{anything}} = 0$$

Suppose my team is good.

$$\text{Win} = \frac{5}{3} = 1.7 \text{ improves } \frac{9}{3} = 3$$

$$\text{More improvement} = \frac{27}{3} = 9$$

that will go to $= \infty$.

$$\text{Odds} = \frac{\text{Something happening}}{\text{Something not happening}} = \frac{P}{1-P}$$

So odds (win) will go from 1 to ∞ .
And odds (loss) will go from 0 to 1.

So, symmetry is not there here which is asymmetry meaning it makes it difficult to compare Odd (win) and Odd (loss).

Suppose, if Odds are against 1 to 6, Odds = $\frac{1}{6} = 0.17$ but in favour, $\frac{6}{1} = 6$

So, logs of Odds will make it symmetry.

$$\text{Log of odd against 1 to 6} = \log(\text{Odds}) = \log\left(\frac{1}{6}\right) = \log(0.17) = -1.79$$

$$\text{Log of odd against 6 to 1} = \log(\text{Odds}) = \log\left(\frac{6}{1}\right) = \log(6) = 1.79$$

In short log make the midpoint to 0.

- There are 2 methods to check Assumption 2.

i) Create a log variable of Independent variable and check p value.

ii) Visualization

1) Box Tidewell Test

- Add log-transform interaction variable between continuous independent variable. and run the model.

- Suppose we have two variables: Age and Fare.

Transform them to natural log: $\log(\text{Age})$ and $\log(\text{Fare})$.

Run the logistic model and check p value.

Age - 0.051

$\log(\text{Age})$ - 0.101

Fare - 0.000

$\log(\text{Fare})$ - 0.000

$\log(\text{Fare}) < 0.001$, statistically significant.

$\log(\text{Age}) = 0.101$, statistically insignificant.

This means there is non-linearity in Fare Feature is < 0.001 , and assumption has been violated.

→ We can resolve this by including a polynomial term (eg Fare^2) to account for non-linearity.

ii) Visual Check

- Check scatter plot between Independent Variable and Log-odds(predicted)

eg, $\text{logit_result} = \text{GLM}(y, X, \text{family} = \text{families.Binomial()}) . \text{fit}()$
 $\text{predicted} = \text{logit_results} . \text{predict}(X)$

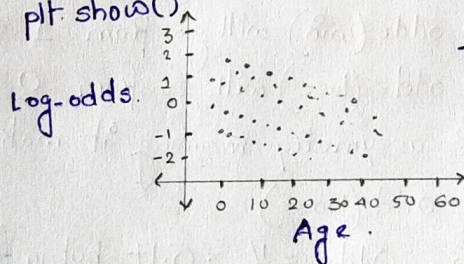
Get log odds value

$\text{log-odds} = \text{np} . \log(\text{predicted} / (1 - \text{predicted}))$

plot for Age variable.

$\text{plt} . \text{scatter}(x = \text{df_titanic}['\text{Age}'] . \text{values}, y = \text{log-odds})$

$\text{plt} . \text{show}()$

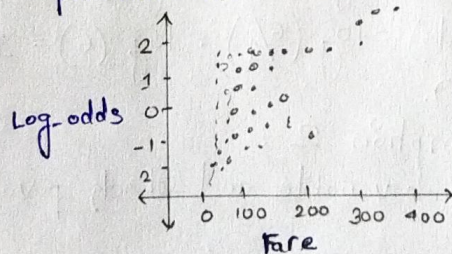


There is a linearity between Age & $\text{log-odds}(\text{predicted})$

plot for Fare variable.

$\text{plt} . \text{scatter}(x = \text{df_titanic}['\text{Fare}'] . \text{values}, y = \text{log-odds})$

$\text{plt} . \text{show}()$



There is non linearity between log-odds and Fare.

So, we will use do $\log(\text{Fare})$ in and check linearity.

ASSUMPTION 3 - NO STRONG INFLUENTIAL OUTLIERS

- It distort accuracy of the model.

Cook Distance \rightarrow

- Cook's distance is the scale change in fitted values, which is useful for identifying in the X values.
- Cook's distance shows the influence of each observations on the fitted response value.
- An observations with Cook distance larger than three/four times the mean Cook's distance might be an outlier.

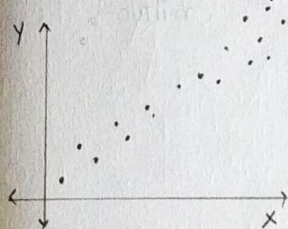
Definition \rightarrow

- Each element in the Cook distance D is the normalized change in fitted response value due to deletion of an observation.

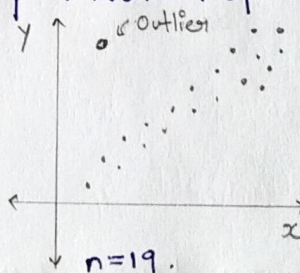
$$D_i = \frac{\sum_{j=1}^n (\hat{y}_j - \hat{y}_{j(i)})^2}{p \text{ MSE}}$$

$\hat{y}_j \rightarrow j^{\text{th}}$ fitted response value, $\hat{y}_{j(i)} \rightarrow j^{\text{th}}$ fitted response value, where the fit does not include observation j .

MSE \rightarrow Mean Square Error, $p \rightarrow$ Number of coefficient in the regression model



Add an outlier \rightarrow



$n=20$;

Coefficient (b_0) = 1.732 [Intercept
Coefficient]

Slope Coefficient (b_1) = 5.117

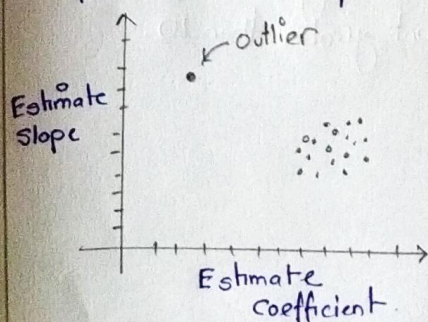
$n=19$.

Intercept Coefficient (b_0) = 8.51

Slope Coefficient (b_1) = 3.32

This estimate change substantially, when there is an outlier.

So what we do is we remove one datapoint at a iteration and check the intercept coefficient and slope coefficient. Then plot the output value of Estimate Slope (b_1) and Estimate Intercept (b_0).



\rightarrow So we use Cook Distance

\rightarrow Cut off is subject to visualization.

\rightarrow Note: Variable have high Cook Distance, does not immediately removed. Check other dimension / angle of the variable

ASSUMPTION 4 - ABSENCE OF MULTICOLLINEARITY.

- Multicollinearity corresponds to a situation where the data contains highly correlated variable.

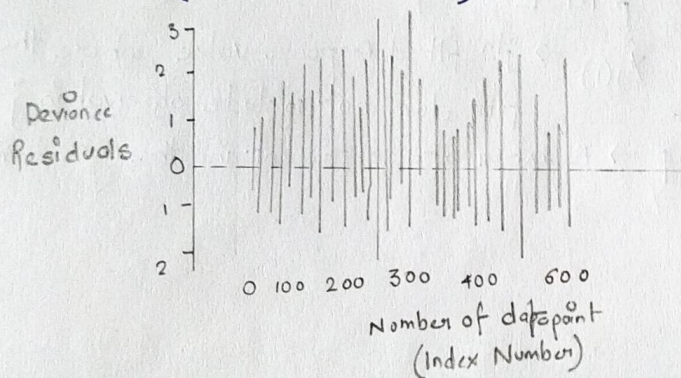
- Leads to reduce the precision of estimated coefficients.

1) **Correlation Matrix** → It checks correlation between independent variables.
- If it exists between combination of more than two independent variables, use VIF.

2) **Variation Inflation Factor** → VIF is equal to the ratio of the overall model variance to the variance of a model that includes only one single independent variable.

ASSUMPTION 5 → INDEPENDENCE OF OBSERVATIONS (AUTOCORRELATION)

- We can also create Residual plot where we plot the deviance residual (Residual Deviance).



ASSUMPTION 6 → SUFFICIENTLY LARGE SAMPLE

- At least 10 observations of each independent category.

Target Column → Survived

Independent Column → Gender (M, F), Travel (Alone, Family)

~~Survived~~ Survived - Yes Survived - No

M 22 11

F 36 2

Alone 13 17

Family 15 16

So all categories (sub) have count greater than 10.