

Legend	
✓ Week 2:	
✓ Week 3:	
✓ Week 4:	
✓ Week 5:	
✓ Week 6:	
✓ Week 7:	
✓ Week 8:	

Load the "song_data.csv" dataset into R.

To validate our model, create two independent dataset, a training dataset & a test dataset, by 50/50 split from the randomly sampling 18,845 observations.

In the training dataset, fit a model based on the scatterplots in EDA, including the predictors, danceability, acousticness, energy, loudness, audio_valence, and tempo, which we suspect a statistical relationship exists among them, with our response, 'song_popularity'.

Check against the two additional conditions with y vs. \hat{y} plot and pairwise scatterplot, which allows us to know which assumption is being violated from the residual plot which will be constructed in the next step.

Refit the model

Check against the four assumptions for linear regression model, which are linearity, uncorrelatedness, homoscedasticity, and normality, by applying the residual plot and the normal Q-Q plot.

Linearity: check if there seems to be a linear association between the 'song_popularity' and the predictors by the residual plot.

Uncorrelatedness: check if the residuals are correlated, where an obvious cluster can be detected on the plot by the residual plot.

Homoscedasticity: check if the residuals have constant variance by the residual plot. (i.e. there is a funnel-like pattern detected)

Normality: Check if the residuals are normally distributed by the Q-Q plot.

if the assumption(s) hold?

Yes

No

Identify the leverage pts or outliers or influential pts by quantifying the amount of influence each observation has by comparing its h_i , r_i , D_i , $DFBETAS_i$, $DFBETAS_i$ against the corresponding cutoff.

Apply the Delta Method for the violation of non-constant variance or/and apply the Box-Cox transformation for the violation of non-linearity or non-normality

if the data error is legitimate

Yes

No

Refit

