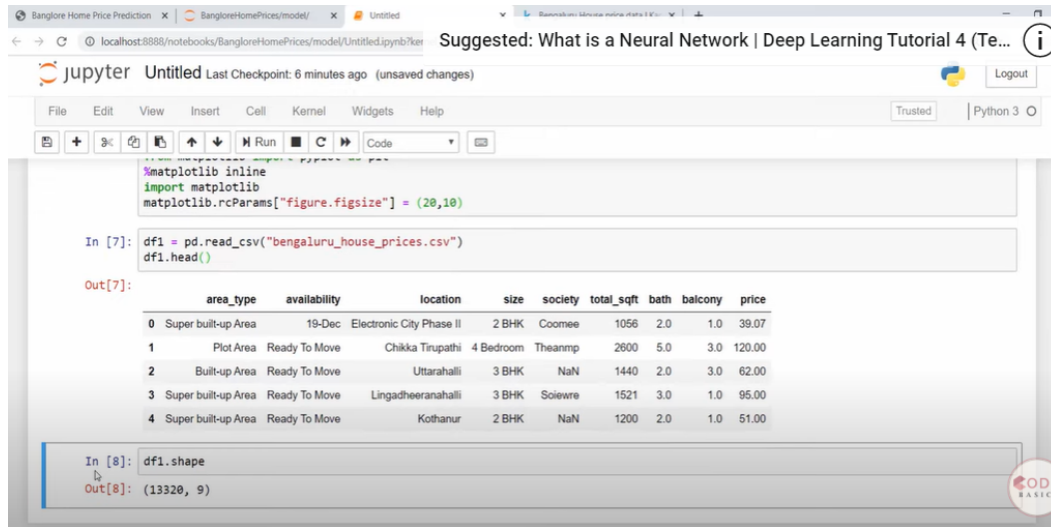We found this cool video on youtube where he is doing data cleaning using different techniques on his real Estate Price Prediction Project. The following link can take you to that video. It shows the model built by sklearn and linear regression. The pandas data frames are loaded and then handled.

https://www.youtube.com/watch?v=_drqJ9SFCgU

Some pictures from the videos are shown below which shows the different techniques used.

Different high level visualization libraries were also used. Some of the interactive graphs were plotted using it.