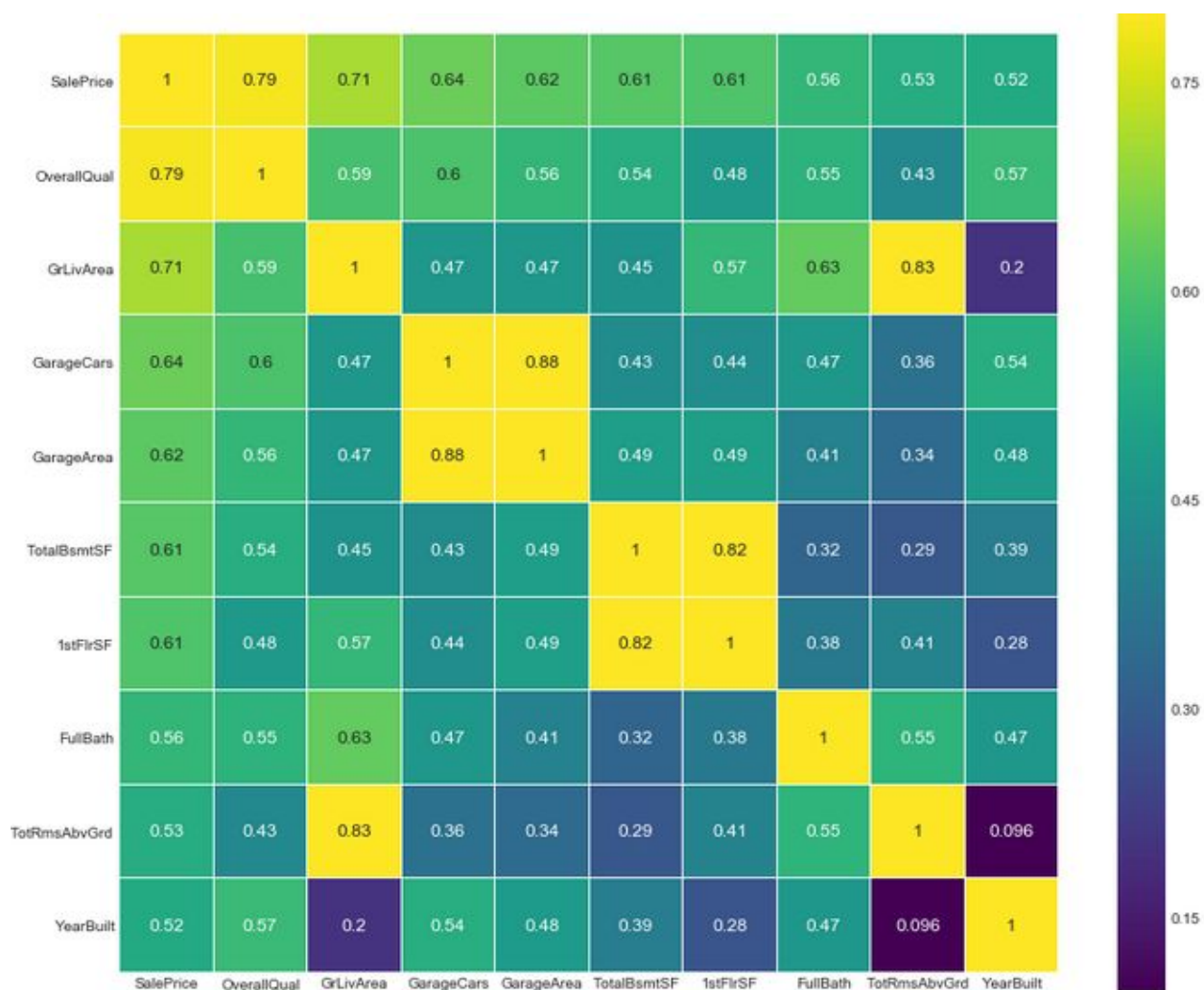I personally think that how this eda and data processing is related to the models that are used to run. RMB that these models are Ridge, Lasso, Kernel Ridge, Elastic net, Gradient Boosting, XGBOost, and LightGBM
These are the questions i have for the script

Dropping out numerical features from looking at correlation
- I am not sure whether or not to drop the features apart from the top 10 correlated with SalesPrice
- So anything that has a correlation less than 0.52 (YearBuilt) can be neglected when modelling



Applying log transformation to the skewed SalesPrice distribution
- Some of the kernels use log transformation, but some uses log1p which is log(1+x) but when I tried both of them it came out identical. Why is that?

Some of the categoricals are not applied with Label Encoder, why?
- I think its because you only apply the function to only ones that are ranks (not fully categorical)
- For example buildingType is not label encoded but extracted as a dummy variable with pd.get_dummies. On the other hand, bsmtCond is label encoded since it has to do with some rank

Filling out missing data
- I noticed that i used the median to fill out LotFrontage but its the median from all_data (train + test). Is it necessary to fill out median separately?

On the BOX COx transformation
- I still dont understand how to choose lambda, but I think its 0.15 is the best one interms of model prediction
- And is 0.75 skewness considered to be bad enuf for transformation?

Dummy Variable trap
- I see there the categoricals are prone to dummy variable trap and multi collinearity presence
- Do we have to deal with it? Or does the model get rid of them automatically