MY SUGGESTIONS:

Question 1

What is the optimal value of alpha for ridge and lasso regression? What will be the changes in the model if you choose double the value of alpha for both ridge and lasso? What will be the most important predictor variables after the change is implemented?

Answer:
- Optimal value of alpha for ridge is 50 and lasso is 0.001
- The performance of the models takes a hit by very slight percentage
- Important predictors for lasso:
        - Positive: GrLivArea, OverallQual, GarageCars, TotRmsAbvGrd, OverallCond
        - Negative: BsmtExposure_None, MSZoning_RM, FireplaceQu_None, MSSubClass, BsmtFinType1_Unf
- Important predictors for ridge:
        - Positive: OverallQual, TotRmsAbvGrd, GarageCars, GrLivArea, FullBath
        - Negetaive: FireplaceQu_None, KitchenQual_TA, Neighborhood_Gilbert, Neighborhood_Edwards, MSZoning_RM

Question 2

You have determined the optimal value of lambda for ridge and lasso regression during the assignment. Now, which one will you choose to apply and why?

Answer:
Model performance for both models looks similar with slight variations, but I will choose lasso. As we have a lot of features, it helps us reduce them and if we can tune the alpha even more minutely it performs better than ridge. Also we can see that the lasso model is fitting to the train data more generically as the performance drop on test data is minimal.

Question 3

After building the model, you realised that the five most important predictor variables in the lasso model are not available in the incoming data. You will now have to create another model excluding the five most important predictor variables. Which are the five most important predictor variables now?

Answer:
Neighborhood_NridgHt, Neighborhood_StoneBr, Neighborhood_Crawfor, BsmtFullBath, Neighborhood_NoRidge

Question 4

How can you make sure that a model is robust and generalisable? What are the implications of the same for the accuracy of the model and why?

Answer:
The more simple the model is more robust it is and will be able to have more generalised properties. When we reduce the complexity of a model there tends to be a drop in the accuracy as the model does not try to fit exactly to the data present for it's training. Thus when we check the accuracy it takes a hit but if we check the drop in accuracy between train and test data, we can clearly see the gain here. The model performance drops to a negligible percentage.