

Subjective Questions

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:

The optimal value of alpha for ridge and lasso regression as below:

- Ridge Alpha

Best optimum • Ridge Alpha 10.0

When Ridge Alpha value is doubles (20) then coefficient **does not** vary much.

Ridge Alpha set to 10.0	Ridge Alpha set to 20.0
<pre>#set alpha to Lowest 0.01 based on the above analysis hs_df_alpha = 10 hs_df_lasso = Ridge(alpha=hs_df_alpha) hs_df_lasso.fit(hs_df_x_train, hs_df_y_train) hs_df_lasso.coef_</pre>	<pre>#set alpha to Lowest 0.01 based on the above analysis hs_df_alpha = 20 hs_df_lasso = Ridge(alpha=hs_df_alpha) hs_df_lasso.fit(hs_df_x_train, hs_df_y_train) hs_df_lasso.coef_</pre>
<pre>array([-4220.82187158, 1902.6897786 , 3938.63390296, 3037.23354194, 11626.5313594 , 5778.35353506, 7834.92417306, 2577.75505749, 4291.39683974, 8409.38247693, 425.32050652, -1261.80792672, 7770.23582669, -579.97811867, 7837.06520233, 10036.82064543, 0. , 14961.04435257, 274.37260133, -1315.80526765, 1243.52707047, 517.11674534, -4693.96027392, 0. , 5951.3196928, 709.07989543, -496.23139609, 2909.99749422, 964.47265703, 603.88213206, 1830.80351631, 1878.72179936, 531.78463218, 2400.9728321, 201.94537492, -365.88581577, 6395.7453026 , 1935.97290936, 3063.3941306 , -2496.57847541, 2412.01775256, -1915.23992052, 1019.78795324, 7551.45300494, -6617.04677102, 1353.80107339, -3173.82271953, 7172.55925024, -3901.89775829, -3093.49703564, -1079.75063655, 882.99811514, -6881.6176168 , 917.62402122, 4114.54669151, 4147.29210516, -5885.02722607, -4697.91461906, 13556.59655448, -5541.56253497, -5220.36581086, -3264.12425229, -422.21262908, -12055.91773534, -7476.19151085, 1328.80715231, -8187.27542497, 8443.87173619, 17486.780105245, -3434.81673178, -1945.80072043, -2799.66669855, -2682.01228705, 3561.99169026, 18852.11854738, -7961.85095958, -1225.87293106, 1224.2534715, 9661.18773846, 2928.86257256,</pre>	<pre>array([-4013.93043732, 2101.11865466, 3746.89152458, 2820.76634 12768.75381851, 5830.12242126, 7470.54733972, 2634.78294 4840.49136944, 8321.9052982 , 489.28564106, -999.61204 7976.96704931, -665.02489366, 7947.12814639, 9198.74674 0. , 14246.66545209, 425.59766766, -1302.95014 1385.60240286, 772.47199237, -4840.01755812, 0. , 6425.07518516, 1060.95225622, -226.34710837, 2577.21906 1605.56343435, 572.18758085, 1732.60412004, 1920.72599 481.01072911, 2285.77302825, 340.43780603, -228.62412 4376.94985352, 872.61178585, 2330.18802041, -2796.58481 2410.99513255, -1567.39853862, 681.06433146, 5466.03522 -4832.56513715, 601.94395673, -1768.92386693, 6276.32013 -2681.91442072, -1759.03215373, -790.98838121, 870.40425 -4345.53579307, 459.77924762, 2261.23555642, 3522.74300</pre>

- lasso Alpha

Best optimum lasso Alpha 100.0

When Lasso Alpha value is doubles (200) then coefficient values **drastically** changed.

Ridge Alpha set to 100.0	Ridge Alpha set to 200.0
<pre>#set alpha to Lowest 0.01 based on the above analysis hs_df_alpha = 100 hs_df_lasso = Lasso(alpha=hs_df_alpha) hs_df_lasso.fit(hs_df_x_train, hs_df_y_train) hs_df_lasso.coef_</pre>	<pre>In [658]: #set alpha to Lowest 0.01 based on the above analysis hs_df_alpha = 200 hs_df_lasso = Lasso(alpha=hs_df_alpha) hs_df_lasso.fit(hs_df_x_train, hs_df_y_train) hs_df_lasso.coef_ Out[658]: array([-4.41867123e+03, 1.94178808e+03, 3.45913772e+03, 2.86977502e+03, 1.32946104e+04, 5.43003153e+03, 9.06822382e+03, 2.47934496e+03, 3.95255491e+03, 8.06254621e+03, 4.92667962e+02, -0.00000000e+00, 6.64392536e+03, -6.93139085e+01, 5.21652401e+02, 0.00000000e+00, 0.00000000e+00, 2.61653431e+04, 4.40944398e+02, -8.73282034e+02, 2.99326426e+01, -0.00000000e+00, -4.38194179e+03, 0.00000000e+00, 5.70060479e+03, 6.28103468e+02, 0.00000000e+00, 2.10342165e+03, 1.93320780e+03, 1.21251216e+01, 1.53311565e+03, 1.69621038e+03, 1.65215552e+02, 1.90973749e+03, 2.03204082e+02, -0.00000000e+00, 0.00000000e+00, 0.00000000e+00, 0.00000000e+00, -1.73955303e+03, 0.00000000e+00, -0.00000000e+00, 0.00000000e+00, 0.00000000e+00, -2.84706289e+03, 0.00000000e+00, -0.00000000e+00, 6.0765379e+03, -0.00000000e+00, -0.00000000e+00, 0.00000000e+00, 0.00000000e+00, -0.00000000e+00, 0.00000000e+00, 0.00000000e+00, 2.1978007e+03, -0.00000000e+00, -0.00000000e+00, 1.71689300e+04, -0.00000000e+00, -0.00000000e+00, -0.00000000e+00, 0.00000000e+00, -3.48055227e+03, -1.73504468e+03, 0.00000000e+00, -3.36591918e+03, 6.56743340e+03, 2.25336348e+04, -3.90074320e+02, 0.00000000e+00, 0.00000000e+00, -0.00000000e+00, 6.15483757e+03, 2.28596761e+04, -0.00000000e+00,</pre>

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:

r2_score computed using Lasso regression is lower (i.e. 0.8610351617897898) compare to Ridge regression (i.e. 0.9289551150785119) hence I would choose Lasso regression.

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:

```
hs_df_train_cols=hs_df_x_train.columns[hs_df_rfe.support_]
|

hs_df_x_train2=hs_df_x_train[hs_df_train_cols]
hs_df_x_train_cons = smapi.add_constant(hs_df_x_train2)

#create first liner regression model
hs_df_lr_md1 =smapi.OLS(hs_df_y_train,hs_df_x_train_cons)

#fit the model
hs_df_lr_md1=hs_df_lr_md1.fit()

#Print the summary
hs_df_lr_md1.summary()
```

OLS Regression Results

Dep. Variable:	SalePrice	R-squared:	0.797
Model:	OLS	Adj. R-squared:	0.793
Method:	Least Squares	F-statistic:	215.2
Date:	Tue, 14 Jun 2022	Prob (F-statistic):	4.64e-290
Time:	03:24:29	Log-Likelihood:	-10582.
No. Observations:	893	AIC:	2.120e+04
Df Residuals:	876	BIC:	2.128e+04
Df Model:	16		
Covariance Type:	nonrobust		

Columns highlighted in the below table shall be excluded.

	coef	std err	t	P> t	[0.025	0.975]
const	1.911e+05	1.1e+04	17.332	0.000	1.69e+05	2.13e+05
BsmtFinSF1	1.431e+04	985.475	14.519	0.000	1.24e+04	1.62e+04
BsmtFinSF2	-823.6565	1138.967	-0.723	0.470	-3059.080	1411.767
BsmtUnfSF	1875.0789	984.265	1.905	0.057	-56.714	3806.872
TotalBsmtSF	1.696e+04	1561.980	10.859	0.000	1.39e+04	2e+04
1stFlrSF	1.264e+04	1637.211	7.721	0.000	9428.338	1.59e+04
2ndFlrSF	1.384e+04	1019.403	13.576	0.000	1.18e+04	1.58e+04
GrLivArea	2.196e+04	897.665	24.467	0.000	2.02e+04	2.37e+04
RoofStyle_Gable	2.285e+04	7257.858	3.148	0.002	8603.689	3.71e+04
RoofStyle_Gambrel	2.316e+04	1.34e+04	1.726	0.085	-3172.183	4.95e+04
RoofStyle_Hip	3.065e+04	7512.534	4.080	0.000	1.59e+04	4.54e+04
RoofStyle_Mansard	-2556.5676	1.67e+04	-0.153	0.878	-3.53e+04	3.02e+04
RoofStyle_Shed	3.767e+04	3.06e+04	1.232	0.218	-2.23e+04	9.77e+04
RoofMatl_Metal	6.885e+04	3.06e+04	2.253	0.024	8876.114	1.29e+05
RoofMatl_Tar&Grv	1.047e+04	1.65e+04	0.633	0.527	-2.2e+04	4.29e+04
Exterior1st_CBlock	-9727.4452	1.88e+04	-0.518	0.604	-4.66e+04	2.71e+04
Exterior2nd_CBlock	-9727.4452	1.88e+04	-0.518	0.604	-4.66e+04	2.71e+04
ExterQual_Fa	-5.718e+04	1.57e+04	-3.646	0.000	-8.8e+04	-2.64e+04
ExterQual_TA	-4.084e+04	2859.455	-14.281	0.000	-4.64e+04	-3.52e+04
GarageQual_Fa	-1.758e+04	1.19e+04	-1.473	0.141	-4.1e+04	5850.066
GarageQual_TA	-8986.9745	1.01e+04	-0.893	0.372	-2.87e+04	1.08e+04
Omnibus:	238.369	Durbin-Watson:	2.016			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	1777.064			
Skew:	1.005	Prob(JB):	0.00			
Kurtosis:	9.612	Cond. No.	2.45e+17			

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The smallest eigenvalue is 4.58e-32. This might indicate that there are strong multicollinearity problems or that the design matrix is singular.

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:

Simple models with test accuracy is not lesser than the training score is the best models. Null values, outliers shall be examined properly as they impact the assessment. Unwanted null values and outlier records shall be removed or imputed. Data Dictionary shall be closely followed to fill null values.

If required, Bias-Variance trade-off shall be used.