

IEOR 165 PROJECT REPORT

TEAM:

Sumaanyu Maheshwari, Ebru Kasikaralar, Shannon Nakamura, Yamac Karakus, Ferdinand Calisir

Part 1: Wine Quality Prediction

Model Background

In this project, we will be using various regression methods to identify the coefficients of a linear model relating wine quality to different features of the wine: fixed acidity, volatile, acidity citric acid, residual sugar, chlorides, free sulfur dioxide, total sulfur dioxide, density, pH, sulphates and alcohol. Our design matrix (X) is composed of the wine features and our response variable (y) will be the quality of wine. We will be using Pandas and ML libraries (Scikitlearn) to build our models. Our models will be based on the methods of ordinary least square, ridge regression, lasso regression, and elastic net. Moreover, we will be including the plots of tuning parameter versus cross-validation error, coefficients labeled by parameter and the minimum squared error or cross-validation error based on the choice of optimal parameters. Throughout our explanation, we will also be attaching our source code. Data from: P. Cortez, A. Cerdeira, F. Almeida, T. Matos, and J. Reis, Modeling wine preferences by data mining from physicochemical properties, Decision Support Systems, vol. 47, no.4:547-553, 2009.

```
In [95]: #importing relevant libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
from pandas.plotting import scatter_matrix
from sklearn import preprocessing
import sklearn.linear_model as lm
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import Ridge
from sklearn.linear_model import Lasso
from sklearn.linear_model import ElasticNet
from sklearn import svm, grid_search
import warnings
warnings.filterwarnings('ignore')
```

In []:

Exploratory Data Analysis

```
In [96]: df = pd.read_csv("winequality-red.csv", sep = ";")
```

```
In [97]: df.head()
```

Out[97]:

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol
0	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4
1	7.8	0.88	0.00	2.6	0.098	25.0	67.0	0.9968	3.20	0.68	9.8
2	7.8	0.76	0.04	2.3	0.092	15.0	54.0	0.9970	3.26	0.65	9.8
3	11.2	0.28	0.56	1.9	0.075	17.0	60.0	0.9980	3.16	0.58	9.8
4	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4

We first checked whether there are any Null values in the data frame. If there were any Null values in the data frame, we would first have to handle them without changing the structure of the data. We realized there is no Null value in the data frame.

```
In [98]: print("Percentage of Null values in each column: ")  
display((df.isnull().sum()/len(df))*100)
```

Percentage of Null values in each column:

```
fixed acidity      0.0  
volatile acidity   0.0  
citric acid        0.0  
residual sugar     0.0  
chlorides          0.0  
free sulfur dioxide 0.0  
total sulfur dioxide 0.0  
density            0.0  
pH                 0.0  
sulphates          0.0  
alcohol            0.0  
quality            0.0  
dtype: float64
```

We also wanted to check the statistics of the features. Following command has given us the mean, min, max, median, 25th percentile and 75th percentile and standard variation of the data regarding each feature.

In [99]: df.describe()

Out[99]:

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide
count	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000
mean	8.319637	0.527821	0.270976	2.538806	0.087467	15.874922	46.467100
std	1.741096	0.179060	0.194801	1.409928	0.047065	10.460157	32.895000
min	4.600000	0.120000	0.000000	0.900000	0.012000	1.000000	6.000000
25%	7.100000	0.390000	0.090000	1.900000	0.070000	7.000000	22.000000
50%	7.900000	0.520000	0.260000	2.200000	0.079000	14.000000	38.000000
75%	9.200000	0.640000	0.420000	2.600000	0.090000	21.000000	62.000000
max	15.900000	1.580000	1.000000	15.500000	0.611000	72.000000	289.000000

Following graphs are the density plots showing the distribution of all the features:

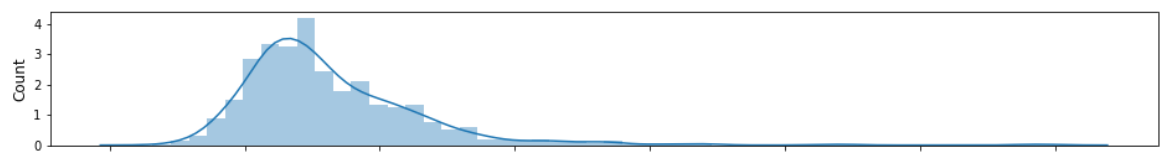
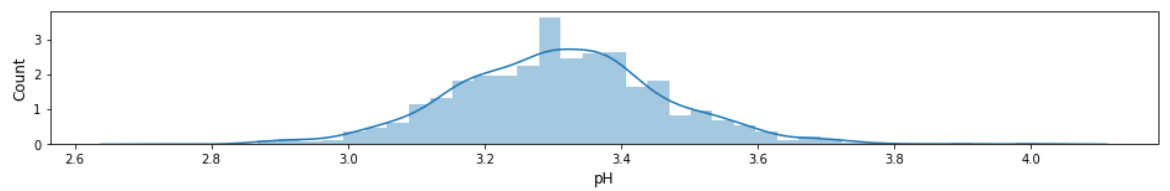
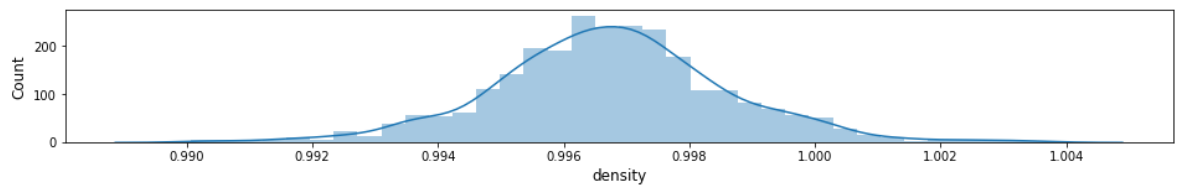
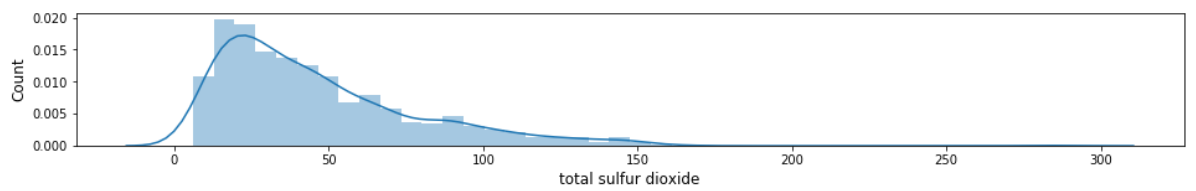
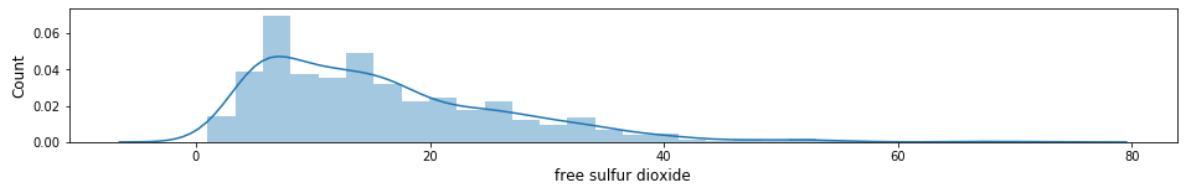
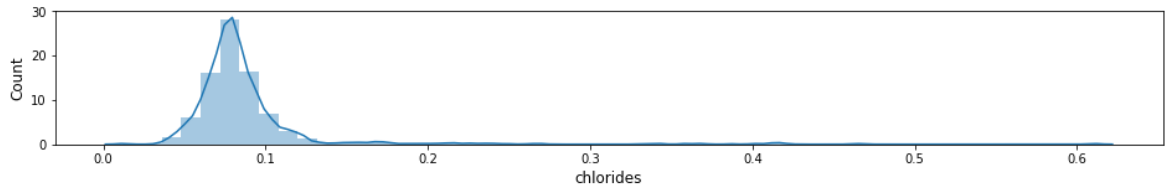
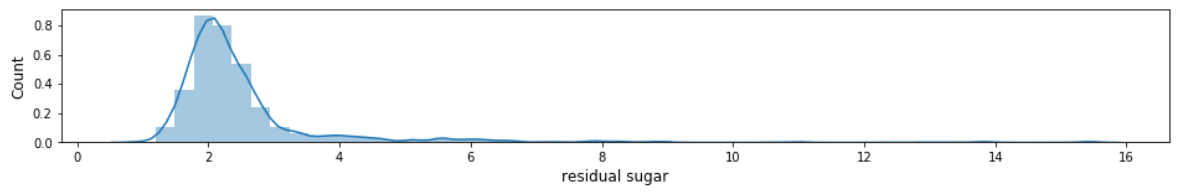
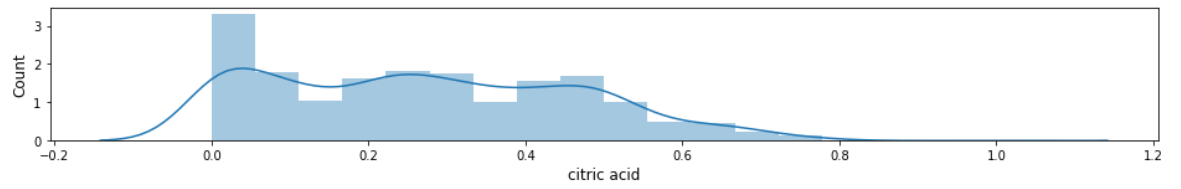
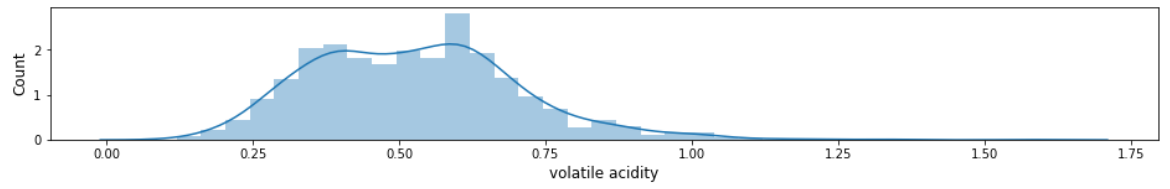
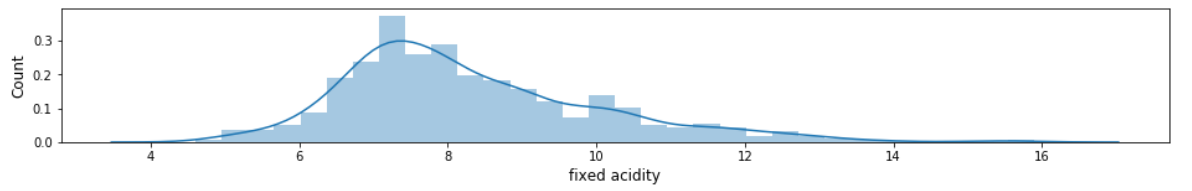
```
In [115]: # check for skewness
new_df=df.drop(columns = "quality")
col_names = new_df.columns

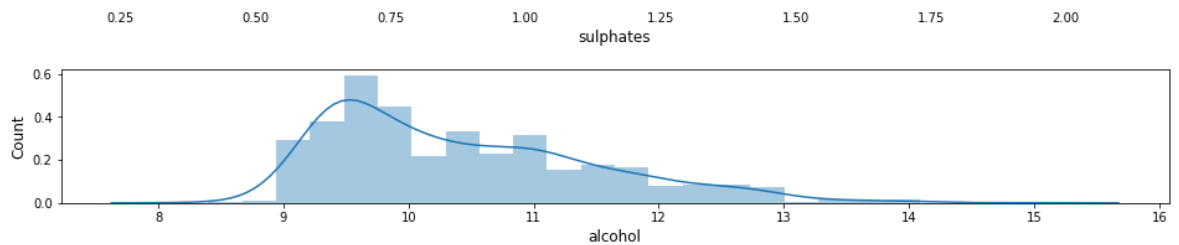
fig, ax = plt.subplots(len(col_names), figsize=(16,12))

for i, col_val in enumerate(col_names):

    sns.distplot(df[col_val], hist=True, ax=ax[i])
    ax[i].set_xlabel(col_val, fontsize=12)
    ax[i].set_ylabel('Count', fontsize=12)

plt.subplots_adjust(bottom = 0.01, top = 2, hspace = 0.5 );
plt.show();
```





From the distributions above, we can see that each feature is measured using a different scale. Therefore, the means and standard deviations of each feature varies significantly from each other. Following data frame tells us the correlations between the features of wine. We wanted to look at the correlation in order to understand whether there is any multicollinearity between some features of the models.

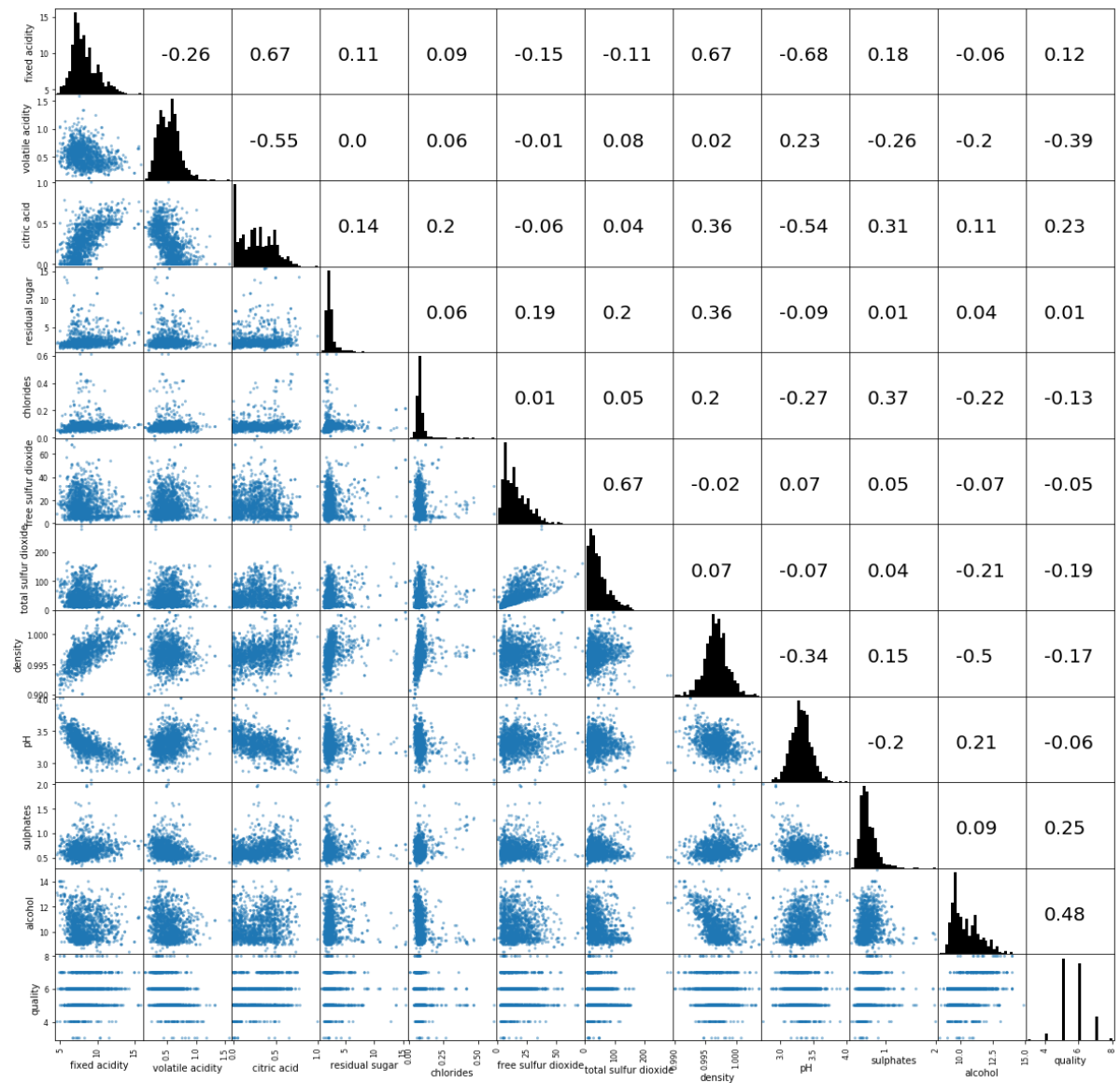
In [101]: `df.corr()`

Out[101]:

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density
fixed acidity	1.000000	-0.256131	0.671703	0.114777	0.093705	-0.153794	-0.113181	0.668047
volatile acidity	-0.256131	1.000000	-0.552496	0.001918	0.061298	-0.010504	0.076470	0.022026
citric acid	0.671703	-0.552496	1.000000	0.143577	0.203823	-0.060978	0.035533	0.364947
residual sugar	0.114777	0.001918	0.143577	1.000000	0.055610	0.187049	0.203028	0.355283
chlorides	0.093705	0.061298	0.203823	0.055610	1.000000	0.005562	0.047400	0.200632
free sulfur dioxide	-0.153794	-0.010504	-0.060978	0.187049	0.005562	1.000000	0.667666	-0.021946
total sulfur dioxide	-0.113181	0.076470	0.035533	0.203028	0.047400	0.667666	1.000000	0.071269
density	0.668047	0.022026	0.364947	0.355283	0.200632	-0.021946	0.071269	1.000000
pH	-0.682978	0.234937	-0.541904	-0.085652	-0.265026	0.070377	-0.066495	-0.341699
sulphates	0.183006	-0.260987	0.312770	0.005527	0.371260	0.051658	0.042947	0.148506
alcohol	-0.061668	-0.202288	0.109903	0.042075	-0.221141	-0.069408	-0.205654	-0.496180
quality	0.124052	-0.390558	0.226373	0.013732	-0.128907	-0.050656	-0.185100	-0.174919

```
In [102]: labs = df.columns
axarr = scatter_matrix(df[labs],figsize=(20,20),hist_kws={'bins':30,'color':'black'}) # create scatter plots and histograms

# the plots in the upper-right corner are repetitive, so replace this with correlations for quicker analysis
for i,axlist in enumerate(axarr):
    for j,ax in enumerate(axlist):
        if i<j:
            ax.cla() # clears an Axes object
            ax.text(x=0.8*ax.get_xlim()[0]+0.2*ax.get_xlim()[1],y=0.6*ax.get_ylim()[0]+0.4*ax.get_ylim()[1],
                    s=np.round(np.corrcoef(df[labs[i]],df[labs[j]])[0,1],2),fontsize=20) # adds text to an Axes object
```



In []:

Models

Ordinary Least Squares (OLS)

Ordinary Least Squares requires us to solve the following optimization problem:

$$\hat{\beta} = \arg \min_{\beta} \|Y - X\beta\|_2^2,$$

We are trying to find the linear model with an intercept, so normally we would append the X data with a column of 1's, but the package we used already does that for us. To set up the OLS model we used the linear model package of Sklearn. We used the mean squared error to calculate the score.

```
In [103]: #setting up the design matrix  
X = df.drop(columns= "quality")  
#setting up the response variable  
y = df["quality"]
```

```
In [104]: #standardizing Xs as the values are measured in a different scale from each other  
standardized_X = preprocessing.scale(X)
```

```
In [105]: #defining a function to measure mean squared error  
def mse(actual_y, predicted_y):  
    return (1/len(actual_y)*sum((actual_y - predicted_y)**2))
```

```
In [106]: from sklearn import linear_model as lm  
#fitting the model  
ols_model = lm.LinearRegression()  
ols_model.fit(standardized_X, y)  
#prediction  
predicted_y = ols_model.predict(standardized_X)
```

```
In [107]: # adjusting output model  
OLS_coefficients_array = np.append(ols_model.coef_, ols_model.intercept_)  
indices = np.append(X.columns,"intercept")  
OLS_coefficients = pd.DataFrame(OLS_coefficients_array,index = indices).rename(column  
ns = {0: "coefficients"})
```



```
In [116]: print("The Coefficients of the Model are as follows:")
          OLS_coefficients
```

The Coefficients of the Model are as follows:

Out[116]:

	coefficients	Ridge	Lasso	Elastic Net
fixed acidity	0.043497	0.053098	0.002904	0.050472
volatile acidity	-0.193967	-0.181537	-0.187431	-0.181117
citric acid	-0.035553	-0.015915	-0.011654	-0.014110
residual sugar	0.023019	0.026914	0.006605	0.025651
chlorides	-0.088183	-0.085085	-0.085485	-0.084713
free sulfur dioxide	0.045606	0.038459	0.035767	0.037505
total sulfur dioxide	-0.107356	-0.101091	-0.098463	-0.100176
density	-0.033737	-0.053216	-0.000000	-0.050883
pH	-0.063842	-0.047636	-0.068041	-0.047515
sulphates	0.155277	0.150007	0.144442	0.149181
alcohol	0.294243	0.269769	0.306785	0.270673
intercept	5.636023	5.636023	5.636023	5.636023

```
In [120]: MSE_linear = mse(y, predicted_y)
          print("MSE of the model is as follows:")
          print("MSE : {}".format(MSE_linear))
```

MSE of the model is as follows:
MSE : 0.4167671672214083.

Ridge Regression

Ridge Regression or L2 - regularization is done by solving the following optimization problem with a parameter λ .

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} ||Y - X\beta||_2^2 + \lambda ||\beta||_2^2$$

In the equation, lambda is the tuning parameter. In order to find the coefficients that minimize the mean squared error, we tried to optimize the value of lambda. First, we built a ridge regression model using Scikit-learn Ridge package. Later on, we created a range of lambdas starting from 0.1 and ending at 100 with increments of 0.05. In order to find the optimal lambda that would result in the minimum mean squared error, we used Scikitlearn. model_selection's GridSearchCV package to generate a 5 fold cross-validation. We also standardized the X values in order to achieve a more accurate model. After building the model, we used .fit function to fit the model on our X and Y values and selecting the optimal lambda based on the result GridSearchCV function returned, we used a fit command to fit the model to our data. The function outputted an array of length 1998 of cross-validation errors for each value of lambda. We then used this to plot cross-validation error vs. value of tuning parameter. This gave us an optimal tuning parameter choice of 84.35. Please refer to the plot below.

```
In [23]: ridge = Ridge()
```

```
In [24]: alphas = np.arange(0.1, 100, 0.05)
```

```
In [25]: parameters = {'alpha': alphas}
```

```
In [26]: r_regressor = GridSearchCV(ridge, parameters, scoring='neg_mean_squared_error', cv = 5)
```

```
In [27]: r_regressor.fit(standardized_X, y)
```

```
Out[27]: GridSearchCV(cv=5, error_score='raise',  
    estimator=Ridge(alpha=1.0, copy_X=True, fit_intercept=True, max_iter=None,  
    normalize=False, random_state=None, solver='auto', tol=0.001),  
    fit_params=None, iid=True, n_jobs=1,  
    param_grid={'alpha': array([ 0.1 , 0.15, ..., 99.9 , 99.95])},  
    pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',  
    scoring='neg_mean_squared_error', verbose=0)
```

```
In [28]: r_regressor.best_score_
```

```
Out[28]: -0.43533034290183703
```

```
In [29]: r_regressor.best_index_
```

```
Out[29]: 1685
```

```
In [30]: alphas[r_regressor.best_index_]
```

```
Out[30]: 84.350000000000002
```

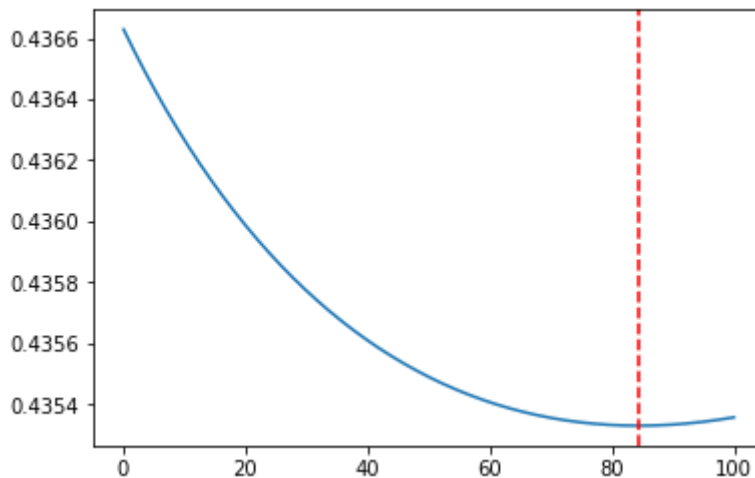
```
In [31]: scores_mean = r_regressor.cv_results_['mean_test_score']
```

```
In [119]: MSE_ridge = min(-1*scores_mean)
print("MSE with optimal lambda:")
print("MSE is {}".format(MSE_ridge))
```

MSE with optimal lambda:
MSE is 0.43533034290183703

```
In [33]: _, ax = plt.subplots(1,1)
ax.plot(alphas, -1*scores_mean)
plt.axvline(x=alphas[r_regressor.best_index_], color='red', linestyle='dashed')
```

Out[33]: <matplotlib.lines.Line2D at 0x7fb399a34d68>



Building the model by inputting optimal alpha:

```
In [34]: ridge_model = Ridge(alpha = alphas[r_regressor.best_index_])
ridge_model.fit(standardized_X,y)
```

Out[34]: Ridge(alpha=84.35000000000002, copy_X=True, fit_intercept=True, max_iter=None, normalize=False, random_state=None, solver='auto', tol=0.001)

```
In [35]: Ridge_coefficients_array = np.append(ridge_model.coef_, ridge_model.intercept_)
indices = np.append(X.columns, "intercept")
Ridge_coefficients = pd.DataFrame(Ridge_coefficients_array, index = indices).rename(columns = {0:"coefficients"})
```

```
In [121]: print("Following are the coefficients of Ridge Regressor:")  
Ridge_coefficients
```

Following are the coefficients of Ridge Regressor:

Out[121]:

	coefficients
fixed acidity	0.053098
volatile acidity	-0.181537
citric acid	-0.015915
residual sugar	0.026914
chlorides	-0.085085
free sulfur dioxide	0.038459
total sulfur dioxide	-0.101091
density	-0.053216
pH	-0.047636
sulphates	0.150007
alcohol	0.269769
intercept	5.636023

Lasso Regression

Lasso Regression or L1-regularization is done by solving the following optimization problem with a parameter λ .

$$\hat{\beta} = \arg \min \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1,$$

Lasso stands for Least Absolute Shrinkage and Selection Operator. It is very similar to the previously described Ridge Regression, except that Lasso penalizes the sum of the absolute weight values instead of the sum of squared values. Intuitively, the exponents in each procedure correspond to the names - L1 and L2 regularization. In the same way as before, lambda is the tuning parameter in this equation. In order to find the coefficients that minimize the mean squared error, we optimized the value of lambda. First, we built a lasso regression model using Scikit-learn Ridge package. Later on, we created a range of lambdas starting from 0.0001 and ending at 0.1 with increments of 0.0005. The rest of the procedure is the same as it was for Ridge -- in order to find the optimal lambda that would result in the minimum mean squared error, we used Scikitlearn. model_selection's GridSearchCV package to generate a 5 fold cross-validation. We also standardized the X values in order to achieve a more accurate model. After building the model, we used .fit function to fit the model on our X and Y values and selecting the optimal lambda based on the result GridSearchCV function returned, we used a fit command to fit the model to our data. The function outputted an array of length 1998 of cross-validation errors for each value of lambda. We then used this to plot cross-validation error vs. value of tuning parameter. This gave us an optimal tuning parameter choice of 0.0047. Please refer to the plot below.

```
In [36]: lasso = Lasso()
```

```
In [37]: alphas = np.arange(0.0001, 0.1, 0.00005)
```

```
In [38]: parameters = {'alpha': alphas}
```

```
In [39]: l_regressor = GridSearchCV(lasso, parameters, scoring='neg_mean_squared_error', cv = 5)
```

```
In [40]: l_regressor.fit(standardized_X, y)
```

```
Out[40]: GridSearchCV(cv=5, error_score='raise',
    estimator=Lasso(alpha=1.0, copy_X=True, fit_intercept=True, max_iter=1000,
    normalize=False, positive=False, precompute=False, random_state=None,
    selection='cyclic', tol=0.0001, warm_start=False),
    fit_params=None, iid=True, n_jobs=1,
    param_grid={'alpha': array([0.0001, 0.00015, ..., 0.0999, 0.09995])},
    pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
    scoring='neg_mean_squared_error', verbose=0)
```

```
In [41]: l_regressor.best_params_
```

```
Out[41]: {'alpha': 0.0047000000000000001}
```

```
In [42]: l_regressor.best_score_
```

```
Out[42]: -0.4356902226391317
```

```
In [43]: alphas[l_regressor.best_index_]
```

```
Out[43]: 0.0047000000000000001
```

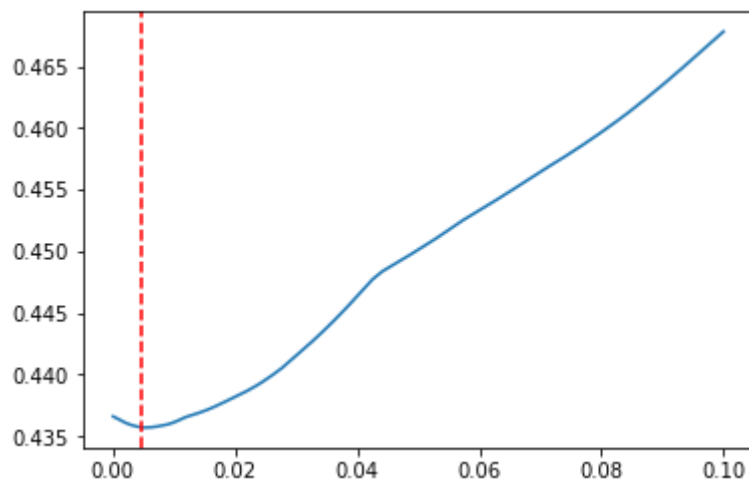
```
In [44]: l_scores_mean = l_regressor.cv_results_['mean_test_score']
```

```
In [122]: MSE_lasso = min(-l_scores_mean)
print("MSE with optimal lambda")
print("MSE is {}".format(MSE_lasso))
```

MSE with optimal lambda
MSE is 0.4356902226391317

```
In [46]: _, ax = plt.subplots(1,1)
ax.plot(alphas, -1*l_scores_mean)
plt.axvline(x=alphas[92], color='red', linestyle='dashed')
```

Out[46]: <matplotlib.lines.Line2D at 0x7fb391f14cc0>



```
In [110]: lasso_model = Lasso(alpha = alphas[l_regressor.best_index_])
lasso_model.fit(standardized_X, y)
Lasso_coefficient_array = np.append(lasso_model.coef_, lasso_model.intercept_)
indices = np.append(X.columns, "intercept")
Lasso_coefficients = pd.DataFrame(Lasso_coefficient_array, index = indices).rename(columns = {0:"coefficients"})
```

```
In [123]: print("Following are the coefficients of Lasso Regressor:")
Lasso_coefficients
```

Following are the coefficients of Lasso Regressor:

Out[123]:

	coefficients
fixed acidity	0.002904
volatile acidity	-0.187431
citric acid	-0.011654
residual sugar	0.006605
chlorides	-0.085485
free sulfur dioxide	0.035767
total sulfur dioxide	-0.098463
density	-0.000000
pH	-0.068041
sulphates	0.144442
alcohol	0.306785
intercept	5.636023

Elastic Net

For Elastic Nets, we have two parameters that we must optimize for - L1 ratio and the Alpha. We choose 40 evenly spaced alpha values from 10^{-10} to 10^{10} . We use ElasticNetCV to cross-validate and get the optimum values.

```
In [48]: en = ElasticNet()
aph = np.logspace(-10,10,40)
L1_vals = [0.01, 0.1, 0.5, 0.7, 0.9, 0.95, 0.99, 1]
en = lm.ElasticNetCV(alphas = aph, l1_ratio = L1_vals, fit_intercept = True, cv = 5)
en.fit(standardized_X, np.ravel(y))
```

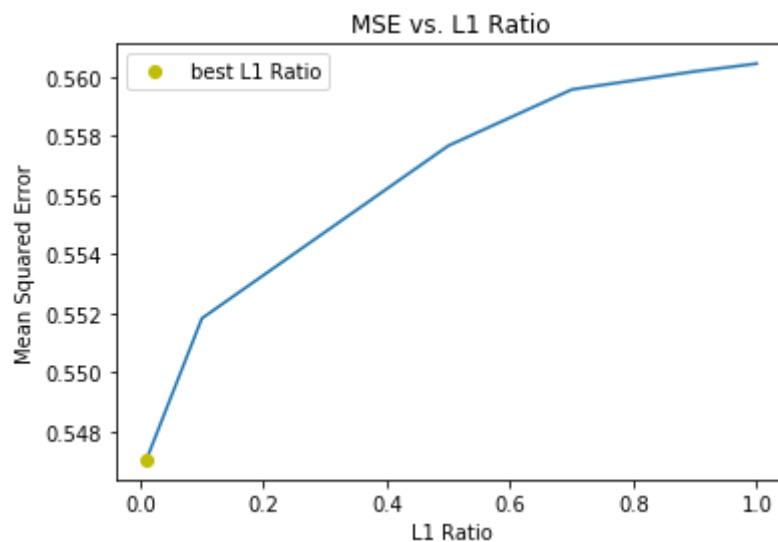
```
Out[48]: ElasticNetCV(alphas=array([1.00000e-10, 3.25702e-10, 1.06082e-09, 3.45511e-09, 1.12534e-08,
3.66524e-08, 1.19378e-07, 3.88816e-07, 1.26638e-06, 4.12463e-06,
1.34340e-05, 4.37548e-05, 1.42510e-04, 4.64159e-04, 1.51178e-03,
4.92388e-03, 1.60372e-02, 5.22335e-02, 1.70125e-01, 5.54102e-01,
..., 8.37678e+06, 2.72833e+07,
8.88624e+07, 2.89427e+08, 9.42668e+08, 3.07029e+09, 1.00000e+10]),
copy_X=True, cv=5, eps=0.001, fit_intercept=True,
l1_ratio=[0.01, 0.1, 0.5, 0.7, 0.9, 0.95, 0.99, 1], max_iter=1000,
n_alphas=100, n_jobs=1, normalize=False, positive=False,
precompute='auto', random_state=None, selection='cyclic',
tol=0.0001, verbose=0)
```

Below is a code snippet and graph for the L1 ratio parameter. It turns out that the L1 ratio is 0.01, meaning our results will be similar to L2 regression.

```
In [49]: l1_mse = []  
for e in en.mse_path_:  
    l1_mse.append(np.mean(e))
```

```
In [50]: plt.plot(L1_vals, l1_mse)  
plt.ylabel("Mean Squared Error")  
plt.xlabel("L1 Ratio")  
plt.title("MSE vs. L1 Ratio")  
plt.plot(en.l1_ratio_, min(l1_mse), "yo", label = "best L1 Ratio")  
plt.legend()
```

Out[50]: <matplotlib.legend.Legend at 0x7fb391f7af28>



```
In [51]: en.l1_ratio_
```

Out[51]: 0.01

```
In [52]: min(l1_mse)
```

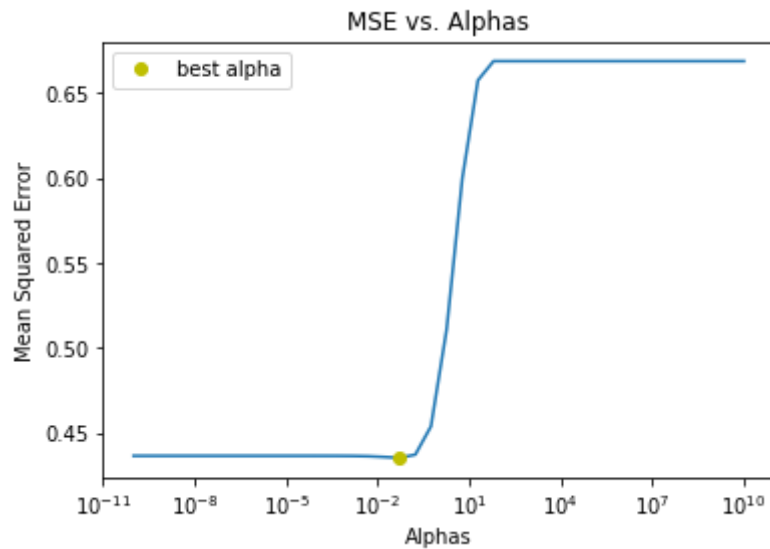
Out[52]: 0.5470344344375717

The MSE we get with this L1 ratio is 0.54. Using the best L1 ratio of 0.01, we now look at the most optimum values for Alpha.

```
In [53]: a_mse = []  
L1_PATH = en.mse_path_[np.argmin(l1_mse)]  
for a in L1_PATH:  
    a_mse.append(np.mean(a))
```



```
In [54]: plt.semilogx(en.alphas_, a_mse)
plt.ylabel("Mean Squared Error")
plt.xlabel("Alphas")
plt.title("MSE vs. Alphas")
plt.plot(en.alpha_, min(a_mse), "yo", label = "best alpha")
plt.legend();
```



```
In [55]: en.alpha_
```

```
Out[55]: 0.05223345074266832
```

```
In [56]: MSE_en = min(a_mse)
print("MSE error is {}".format(MSE_en))
```

```
MSE error is 0.43544100937156677
```

We get alpha as 0.05 and L1 ratio as 0.01. Using these parameters we get the lowest MSE of 0.435

```
In [57]: en_coefficients_array = np.append(en.coef_, en.intercept_)
indices = np.append(X.columns, "intercept")
en_coefficients = pd.DataFrame(en_coefficients_array, index = indices).rename(columns =
{0:"coefficients"})
```

Now, we look at the coefficient terms and the intercept that the model with the abovementioned parameter values provide. Following are the coefficients of Elastic Net:

```
In [124]: print("The coefficients of Elastic Net:")
          en_coefficients
```

The coefficients of Elastic Net:

Out[124]:

	coefficients
fixed acidity	0.050472
volatile acidity	-0.181117
citric acid	-0.014110
residual sugar	0.025651
chlorides	-0.084713
free sulfur dioxide	0.037505
total sulfur dioxide	-0.100176
density	-0.050883
pH	-0.047515
sulphates	0.149181
alcohol	0.270673
intercept	5.636023

Final Comparisons

```
In [112]: OLS = OLS_coefficients
          coefficients = pd.DataFrame(OLS)
          coefficients["Ridge"] = Ridge_coefficients
          coefficients["Lasso"] = Lasso_coefficients
          coefficients["Elastic Net"] = en_coefficients
          coefficients = coefficients.rename(columns = {"coefficients": "OLS"})
```

```
In [125]: print("Coefficients of each model")
coefficients
```

Coefficients of each model

Out[125]:

	OLS	Ridge	Lasso	Elastic Net
fixed acidity	0.043497	0.053098	0.002904	0.050472
volatile acidity	-0.193967	-0.181537	-0.187431	-0.181117
citric acid	-0.035553	-0.015915	-0.011654	-0.014110
residual sugar	0.023019	0.026914	0.006605	0.025651
chlorides	-0.088183	-0.085085	-0.085485	-0.084713
free sulfur dioxide	0.045606	0.038459	0.035767	0.037505
total sulfur dioxide	-0.107356	-0.101091	-0.098463	-0.100176
density	-0.033737	-0.053216	-0.000000	-0.050883
pH	-0.063842	-0.047636	-0.068041	-0.047515
sulphates	0.155277	0.150007	0.144442	0.149181
alcohol	0.294243	0.269769	0.306785	0.270673
intercept	5.636023	5.636023	5.636023	5.636023

```
In [60]: MSEs = [MSE_linear, MSE_ridge, MSE_lasso, MSE_en]
MSE_TABLE = pd.DataFrame(MSEs, index=["OLS", "Ridge", "Lasso", "Elastic Net"]).rename(columns={0:"MSE"})
```

```
In [126]: print("MSE of each model")
MSE_TABLE
```

MSE of each model

Out[126]:

	MSE
OLS	0.416767
Ridge	0.435330
Lasso	0.435690
Elastic Net	0.435441

In []:

```
In [62]: #importing the relevant libraries
import sklearn as skl
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline
import numpy as np
import sklearn.model_selection
from sklearn.model_selection import cross_val_score
import sklearn.svm as svm
from sklearn import preprocessing
from sklearn.metrics import mean_squared_error as mse
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
import pylab as pl
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
```

Part 2: Sales Channel Classification

Model Background

The data for this chapter was obtained from the publication, N. Abreu, Analise do perl do cliente Recheio e Desenvolvimento de um sistema promocional, Mestrado em Marketing, ISCTE-IUL, Lisbon, 2011. The data shows the sales amount of a variety of products, the channel and the region of each retailer. We are constructing various SVMs with different types of kernels in order to predict the sales channel. In this project we are applying three different methods: linear SVM, SVM with Polynomial kernel, and SVM with Gaussian kernel in order to classify data. We are cross-validating parameters for each SVM so that we choose the best parameters that would return models that maximize accuracy or in other words, minimizes cross-validation error.

Exploratory Data Analysis

```
In [63]: df_svm = pd.read_csv('wholesale-customers.csv', sep=',')
df_svm.head()
```

Out[63]:

	Channel	Region	Fresh	Milk	Grocery	Frozen	Detergents_Paper	Delicassen
0	2	3	12669	9656	7561	214	2674	1338
1	2	3	7057	9810	9568	1762	3293	1776
2	2	3	6353	8808	7684	2405	3516	7844
3	1	3	13265	1196	4221	6404	507	1788
4	2	3	22615	5410	7198	3915	1777	5185

We first checked whether there are any Null values in the data frame. If there were any Null values in the data frame, we would first have to handle them without changing the structure of the data. We realized there is no Null value in the data frame.

```
In [64]: print("Percentage of Null values in each column: ")  
display((df_svm.isnull().sum()/len(df_svm))*100)
```

Percentage of Null values in each column:

```
Channel      0.0  
Region       0.0  
Fresh        0.0  
Milk         0.0  
Grocery      0.0  
Frozen       0.0  
Detergents_Paper  0.0  
Delicassen   0.0  
dtype: float64
```

Creating DataFrame for Predictor and Response Variables:

```
In [65]: X = df_svm.iloc[:,1:]  
y = df_svm.iloc[:,1]
```

```
In [66]: X.head()
```

Out[66]:

	Region	Fresh	Milk	Grocery	Frozen	Detergents_Paper	Delicassen
0	3	12669	9656	7561	214	2674	1338
1	3	7057	9810	9568	1762	3293	1776
2	3	6353	8808	7684	2405	3516	7844
3	3	13265	1196	4221	6404	507	1788
4	3	22615	5410	7198	3915	1777	5185

```
In [67]: y.head()
```

Out[67]:

	Channel
0	2
1	2
2	2
3	1
4	2

Processing the data to fit the SVM paradigm. This means that the labels for the Response variable must be between +1 and -1. The math works out much more cleanly if we do things thing way. Since our Channel labels are 1 and 2, we must transform them to give +1 and -1. The transformation we apply is $(-1)^{\text{Channel Label}}$. We also standardize our results. When we are running SVMs, the estimated weights will update similarly rather than at different rates during the build process. This will give us more accurate results. Also, the values in the feature/design matrix (X) were standardized using Sklearn.preprocessing.scale package.

```
In [68]: transformed_y = (-1)**y
transformed_y.head()
```

Out[68]:

	Channel
0	1
1	1
2	1
3	-1
4	1

```
In [69]: standardized_X_values = preprocessing.scale(X)
standardized_X = pd.DataFrame(data = standardized_X_values)
standardized_X.columns = X.columns.values
standardized_X.head()
```

Out[69]:

	Region	Fresh	Milk	Grocery	Frozen	Detergents_Paper	Delicassen
0	0.590668	0.052933	0.523568	-0.041115	-0.589367	-0.043569	-0.066339
1	0.590668	-0.391302	0.544458	0.170318	-0.270136	0.086407	0.089151
2	0.590668	-0.447029	0.408538	-0.028157	-0.137536	0.133232	2.243293
3	0.590668	0.100111	-0.624020	-0.392977	0.687144	-0.498588	0.093411
4	0.590668	0.840239	-0.052396	-0.079356	0.173859	-0.231918	1.299347

Linear SVM

For Linear SVM, we have to optimize for one tuning parameter - the penalty parameter for the error term, C. We set a range of C values by using `np.logspace(-10,0,10)` which returns numbers starting from 10^{-10} and stopping at 10 spaced evenly on a log scale. We then employ the 10 fold cross validation over different values of C to separately determine the optimum parameter, which maximizes the accuracy or minimizes the cross-validation error.

```
In [70]: svc = svm.SVC(kernel='linear')
C_s = np.logspace(-10, 0, 10)
scores = list()
for C in C_s:
    svc.C = C
    cv_scores = cross_val_score(svc, standardized_X, np.ravel(y), cv=10, n_jobs=-1)
    scores.append(np.mean(cv_scores))
```

```
In [71]: scores
```

```
Out[71]: [0.677307963354475,
0.677307963354475,
0.677307963354475,
0.677307963354475,
0.677307963354475,
0.677307963354475,
0.6795806906272024,
0.8318287526427062,
0.8999741602067184,
0.9022515856236787]
```

```
In [72]: errors = []
for i in range(len(scores)):
    errors.append(1-scores[i])
error_linear = errors
error_linear
```

```
Out[72]: [0.32269203664552504,
0.32269203664552504,
0.32269203664552504,
0.32269203664552504,
0.32269203664552504,
0.32269203664552504,
0.3204193093727976,
0.16817124735729383,
0.10002583979328161,
0.09774841437632131]
```

The following command gives the optimum C which maximizes the accuracy or minimum accuracy error of the model. The maximum accuracy is 0.90225158 and minimum accuracy error is 0.09774842.

```
In [73]: best_c = C_s[np.argmax(scores)]
best_c
```

```
Out[73]: 1.0
```

Minimum misclassification error

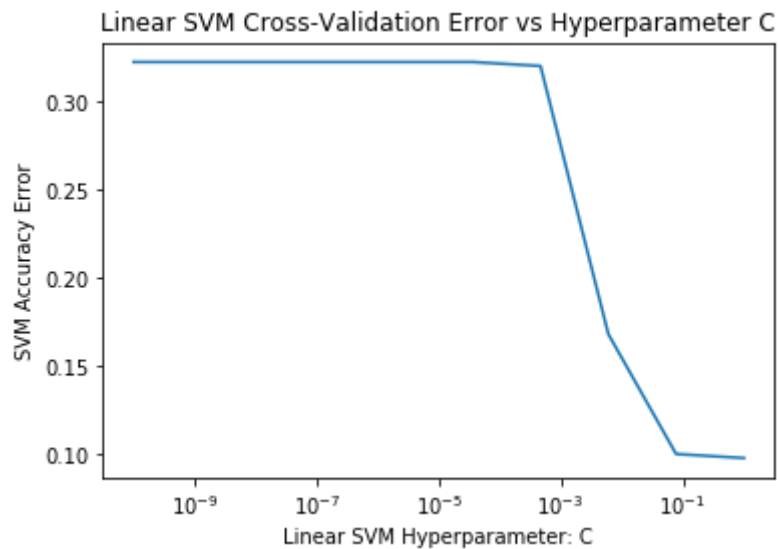
```
In [74]: print("Cross_validation error: {}".format(min(error_linear)))
```

```
Cross_validation error: 0.09774841437632131
```

From above, we can see that the optimal C value is 1. Also, the following plot shows how setting different values of C changes the cross-validation error.

```
In [75]: _, ax = plt.subplots(1,1)
ax.semilogx(C_s, errors)
plt.ylabel("SVM Accuracy Error")
plt.xlabel("Linear SVM Hyperparameter: C")
plt.title("Linear SVM Cross-Validation Error vs Hyperparameter C")
```

Out[75]: Text(0.5, 1.0, 'Linear SVM Cross-Validation Error vs Hyperparameter C')



Using the optimal value of C, which is 1, we set the model and use .fit command to fit the model to our data. The coefficients and the intercept of the model that has the minimum accuracy error are as follows:

```
In [76]: #fitting the model to our data
l_svm = svm.SVC(kernel='linear', C= 1)
model = l_svm.fit(standardized_X, np.ravel(y))
```



```
In [77]: linearsvm_coefficients_array = np.append(model.coef_, model.intercept_)
indices = np.append(standardized_X.columns,"intercept")
linearsvm_coefficients = pd.DataFrame(linearsvm_coefficients_array,index = indices).rename(columns = {0: "coefficients"})
linearsvm_coefficients
```

Out[77]:

	coefficients
Region	0.236901
Fresh	0.035426
Milk	0.346780
Grocery	0.512052
Frozen	-0.281653
Detergents_Paper	2.556616
Delicassen	-0.220582
intercept	0.023486

SVM with Polynomial Kernel

For Polynomial Kernels, we have to optimize for 2 tuning parameters - the penalty parameter for the error term, C, and the degree of the polynomial kernel function, degree. We then employ the 10 fold cross-validation to separately determine the optimum tuning parameters. We also start with default C = 1. We use np.ravel with transformed_y to get a contiguous flattened array for smoother processing.

```
In [78]: degrees = [0,1,2,3,4,5,6]
```

```
In [79]: scores = []
for d in degrees:
    n_scores = cross_val_score(svm.SVC(kernel='poly', C=1, degree = d), standardized_X, n
p.ravel(transformed_y), cv = 10)
    mean_score = np.mean(n_scores)
    scores.append(mean_score)
```

```
In [80]: scores
```

```
Out[80]: [0.677307963354475,
0.9043680996006577,
0.7661674888419073,
0.7956647874089734,
0.754698144233028,
0.7478799624148462,
0.7456048860700023]
```

We now find the errors and plot them

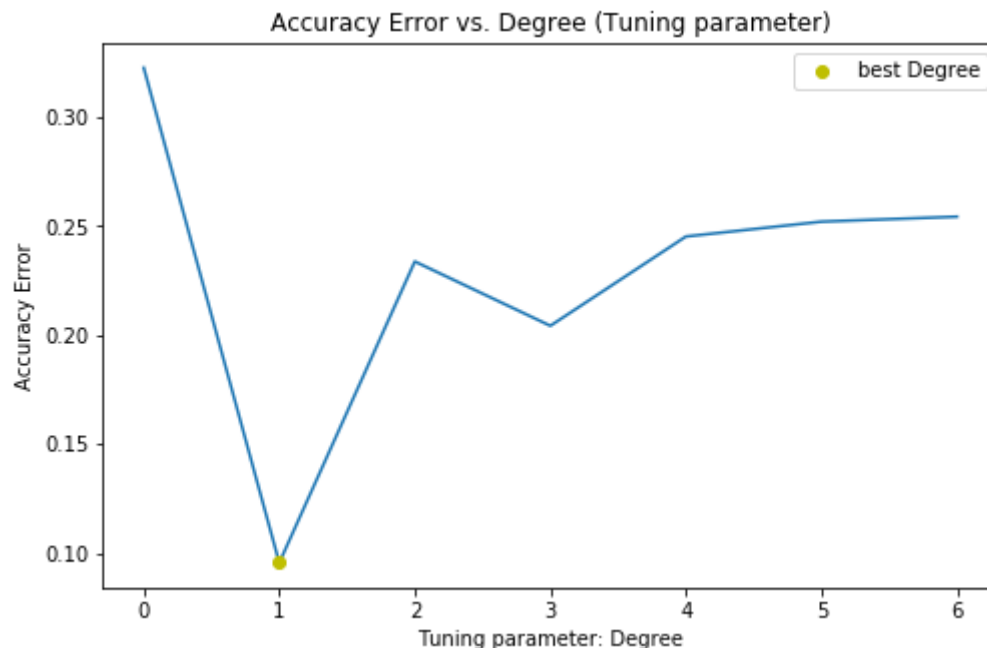
```
In [81]: errors = (np.array(1) - scores)
```

```
In [82]: errors
```

```
Out[82]: array([0.32269204, 0.0956319 , 0.23383251, 0.20433521, 0.24530186,  
              0.25212004, 0.25439511])
```

```
In [83]: best_degree = degrees[np.argmin(errors)]
```

```
In [84]: plt.figure(figsize=(8,5))  
plt.plot(degrees, errors)  
plt.plot(best_degree, min(errors), "yo", label = "best Degree")  
plt.title("Accuracy Error vs. Degree (Tuning parameter)")  
plt.xlabel("Tuning parameter: Degree")  
plt.ylabel("Accuracy Error")  
plt.legend();
```



From the above graph, it is clear the best degree is 1. The cross-validation error associated with degree 1 was 0.09785412. We will use this value and find the most optimum penalty parameter for the error term, C. We choose the most optimum C from $C_s = \text{np.logspace}(-2, 5, 10)$.

```
In [85]: Cs = np.logspace(-2,5,10)  
s = []  
for c in Cs:  
    n_scores = cross_val_score(svm.SVC(kernel='poly', C = c, degree = 1), standardized_X, n  
    p.ravel(transformed_y), cv = 10)  
    mean_n_scores = np.mean(n_scores)  
    s.append(mean_n_scores)
```

```
In [86]: n_errors = (np.array(1) - s )  
n_errors
```

```
Out[86]: array([0.28167489, 0.15003759, 0.10457129, 0.09330632, 0.100074 ,  
0.09552854, 0.09780127, 0.09780127, 0.09552854, 0.09552854])
```

Minimum Misclassification Error

```
In [87]: print("Cross-validation error: {}".format(min(n_errors)))
```

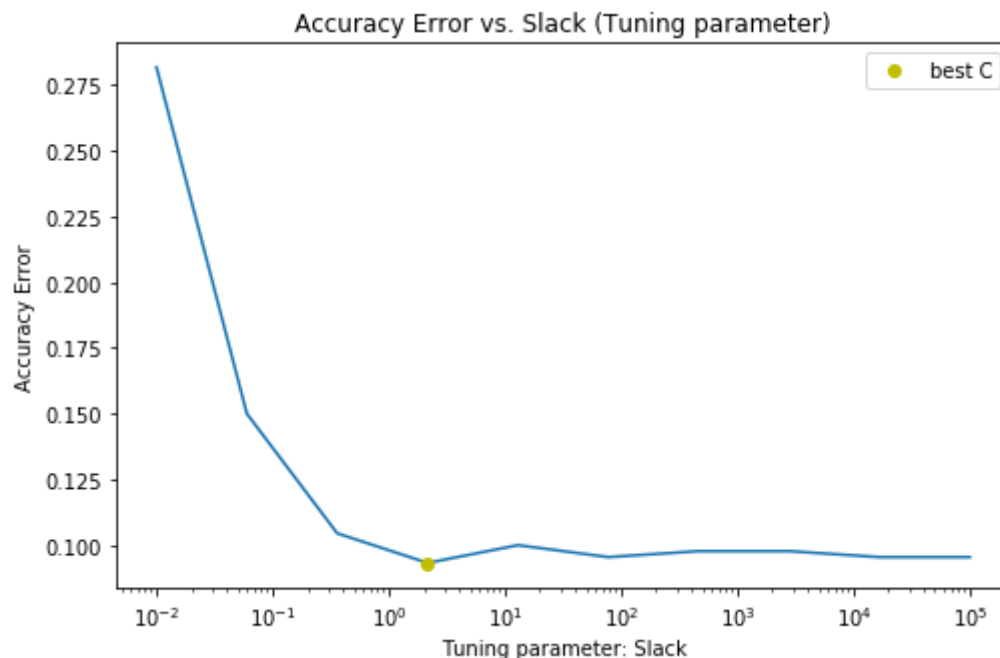
```
Cross-validation error: 0.09330631900399344
```

```
In [88]: best_C = Cs[np.argmin(n_errors)]  
best_C
```

```
Out[88]: 2.1544346900318843
```

We can see that the best C is 2.154 with the lowest cross-validation error of 0.0933. Below is a visual representation of this.

```
In [89]: plt.figure(figsize=(8,5));  
plt.semilogx(Cs, n_errors)  
plt.plot(best_C, min(n_errors), "yo", label = "best C");  
plt.legend()  
plt.title("Accuracy Error vs. Slack (Tuning parameter)")  
plt.xlabel("Tuning parameter: Slack")  
plt.ylabel("Accuracy Error");
```



Results : We found that SVM with degree 1, with default $C = 1$ gave the smallest cross-validation error of 0.09785412. Using this optimized value of degree = 1, we optimized for C . We found that SVM with degree 1, and $C = 2.154$ gave us the lowest cross-validation error of 0.0933.

SVM with Gaussian Kernel

The X data was standardized already in the previous parts. Classifications are 1, and -1 as in the aforementioned processes. SVM with Gaussian kernel has two parameters (C and γ) and these were tuned for minimum accuracy error using cross-validation by performing a grid search. Initially, we performed cross-validation across a larger range of C values, once we had established that C should be in a range of 10^0 to 10^1 we performed another grid search on this interval. The results were displayed as a heatmap as we would multiple lines would have been needed in a 2D graph.

In []:

Final Comparisons

```
In [91]: Tuning_parameter_C = [1, 2.514, 6]
Tuning_Degree = ["-", "1", "-"]
Tuning_Gamma = ["Scale: 1/(n_features*X.var())", "Scale: 1/(n_features*X.var())", 0.01]
Cross_Validation_Error = [0.097748, 0.093306, 0.086363636]
comparison = pd.DataFrame(Tuning_parameter_C, index = ["Linear Kernel", "Polynomial Kernel", "Gaussian Kernel"])
comparison.rename(columns= {0:"Tuning Parameter C"})
comparison["Tuning Parameter Degree"] = Tuning_Degree
comparison["Tuning Parameter Gamma"] = Tuning_Gamma
comparison["Cross Validation Error"] = Cross_Validation_Error
comparison.index.name = "Type of SVM"
comparison
```

Out[91]:

	Tuning Parameter C	Tuning Parameter Degree	Tuning Parameter Gamma	Cross Validation Error
Type of SVM				
Linear Kernel	1.000	-	Scale: 1/(n_features*X.var())	0.097748
Polynomial Kernel	2.514	1	Scale: 1/(n_features*X.var())	0.093306
Gaussian Kernel	6.000	-	0.01	0.086364

In []: