

# Support Vector Machines and Kernels



## ✓ Learning Objectives

- Apply support vector machines (SVMs)—a popular algorithm used for classification problems
- Recognize SVM similarity to logistic regression
- Compute the cost function of SVMs
- Apply regularization in SVMs and some tips to obtain non-linear classifications with SVMs
- Apply Intel® Extension for Scikit-learn\* to leverage underlying compute capabilities of hardware

## scikit-learn\*

Frameworks provide structure that Data Scientists use to build code. Frameworks are more than just libraries, because in addition to callable code, frameworks influence how code is written.

A main virtue of using an optimized framework is that code runs faster. Code that runs faster is just generally more convenient but when we begin looking at applied data science and AI models, we can see more material benefits. Here you will see how optimization, particularly hyperparameter optimization can benefit more than just speed.

These exercises will demonstrate how to apply **the Intel® Extension for Scikit-learn\***, a seamless way to speed up your Scikit-learn application. The acceleration is achieved through the use of the Intel® oneAPI Data Analytics Library (oneDAL). Patching is the term used to extend scikit-learn with Intel optimizations and makes it a well-suited machine learning framework for dealing with real-life problems.

To get optimized versions of many Scikit-learn algorithms using a `patch()` approach consisting of adding these lines of code after importing `sklearn`:

- **`from sklearnex import patch_sklearn`**
- **`patch_sklearn()`**

This exercise relies on installation of Intel® Extension for Scikit-learn\*

If you have not already done so, follow the instructions from Week 1 for instructions

## ✓ Introduction

We will be using the wine quality data set for these exercises. This data set contains various chemical properties of wine, such as acidity, sugar, pH, and alcohol. It also contains a quality metric (3-9, with highest being better) and a color (red or white). The name of the file is `Wine_Quality_Data.csv`.

```
from __future__ import print_function
import os
data_path = [ 'data']
```

```
from sklearn.svm import LinearSVC
from sklearn.svm import SVC
from sklearn.kernel_approximation import Nystroem
from sklearn.svm import SVC
from sklearn.linear_model import SGDClassifier
from sklearn.preprocessing import MinMaxScaler
```

```
from sklearnex import patch_sklearn
patch_sklearn()
```

Intel(R) Extension for Scikit-learn\* enabled (<https://github.com/intel/scikit-learn-intelex>)

## ✓ Question 1

- Import the data.
- Create the target variable `y` as a 1/0 column where 1 means red.
- Create a `pairplot` for the dataset.
- Create a bar plot showing the correlations between each column and `y`
- Pick the most 2 correlated fields (using the absolute value of correlations) and create `x`
- Use `MinMaxScaler` to scale `x`. Note that this will output a `np.array`. Make it a `DataFrame` again and rename the columns appropriately.

```
import pandas as pd
import numpy as np
```

```
filepath = os.sep.join(data_path + ['Wine_Quality_Data.csv'])
data = pd.read_csv(filepath, sep=',')
```

```
y = (data['color'] == 'red').astype(int)
fields = list(data.columns[:-1]) # everything except "color"
correlations = data[fields].corrwith(y)
correlations.sort_values(inplace=True)
correlations
```

```
total_sulfur_dioxide    -0.700357
free_sulfur_dioxide     -0.471644
residual_sugar          -0.348821
citric_acid             -0.187397
```

quality	-0.119323
alcohol	-0.032970
pH	0.329129
density	0.390645
fixed_acidity	0.486740
sulphates	0.487218
chlorides	0.512678
volatile_acidity	0.653036

dtype: float64

```
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

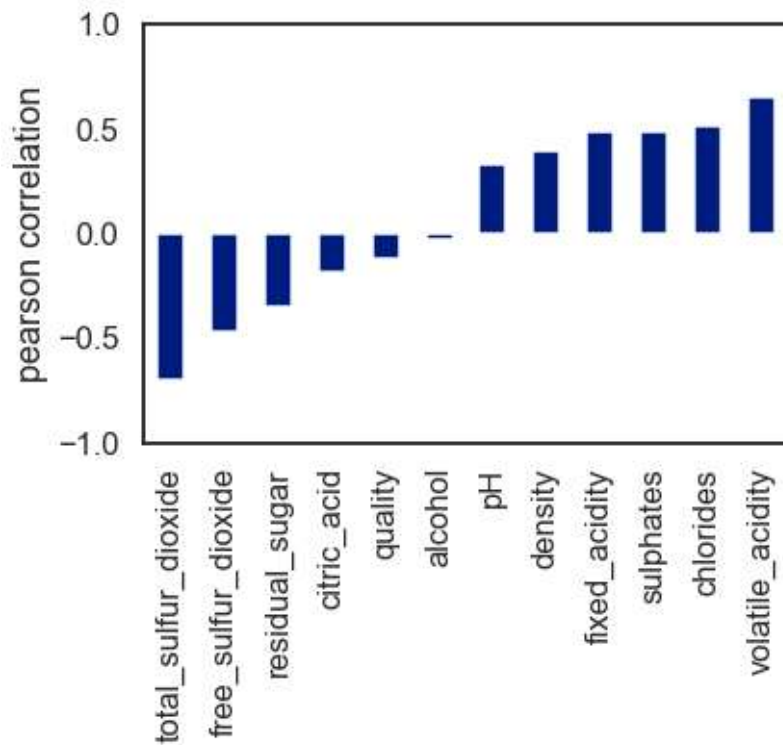
```
sns.set_context('talk')
sns.set_palette('dark')
sns.set_style('white')
```

```
sns.pairplot(data, hue='color')
```

&lt;seaborn.axisgrid.PairGrid at 0x1b3098b5a88&gt;



```
ax = correlations.plot(kind='bar')
ax.set(ylim=[-1, 1], ylabel='pearson correlation');
```



```
fields = correlations.map(abs).sort_values().iloc[-2:].index
print(fields)
X = data[fields]
scaler = MinMaxScaler()
X = scaler.fit_transform(X)
X = pd.DataFrame(X, columns=['%s_scaled' % fld for fld in fields])
print(X.columns)

Index(['volatile_acidity', 'total_sulfur_dioxide'], dtype='object')
Index(['volatile_acidity_scaled', 'total_sulfur_dioxide_scaled'], dtype='object')
```

## ✓ Question 2

The goal for this question is to look at the decision boundary of a LinearSVC classifier on this dataset. Check out [this example](#) in sklearn's documentation.

- Fit a Linear Support Vector Machine Classifier to  $x$ ,  $y$ .
- Pick 300 samples from  $x$ . Get the corresponding  $y$  value. Store them in variables  $x\_color$  and  $y\_color$ . This is because original dataset is too large and it produces a crowded plot.
- Modify  $y\_color$  so that it has the value "red" instead of 1 and 'yellow' instead of 0.
- Scatter plot  $X\_color$ 's columns. Use the keyword argument "color= $y\_color$ " to color code samples.
- Use the code snippet below to plot the decision surface in a color coded way.

```
x_axis, y_axis = np.arange(0, 1, .005), np.arange(0, 1, .005)
xx, yy = np.meshgrid(x_axis, y_axis)
xx_ravel = xx.ravel()
yy_ravel = yy.ravel()
X_grid = pd.DataFrame([xx_ravel, yy_ravel]).T
y_grid_predictions = *[YOUR_MODEL].predict(X_grid)
y_grid_predictions = y_grid_predictions.reshape(xx.shape)
ax.contourf(xx, yy, y_grid_predictions, cmap=plt.cm.autumn_r, alpha=.3)
```

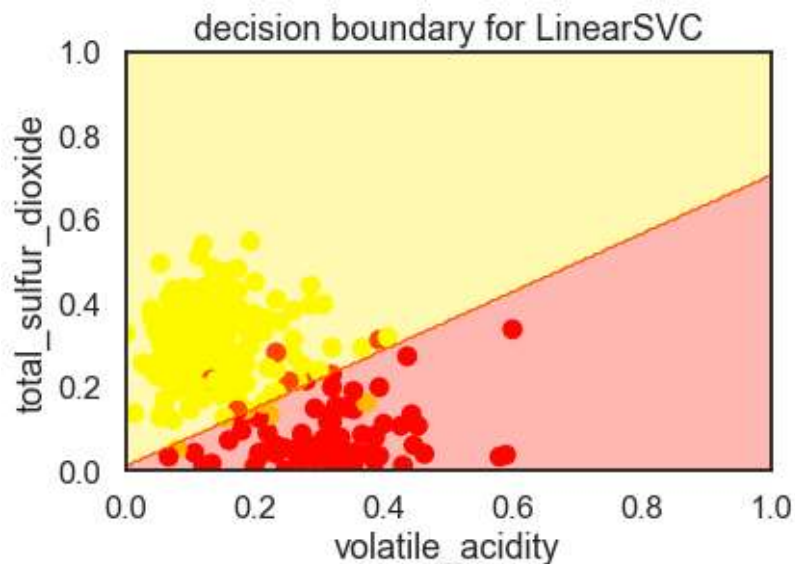
Feel free to experiment with different parameter choices for LinearSVC and see the decision boundary.

```

LSVC = LinearSVC()
LSVC.fit(X, y)

X_color = X.sample(300, random_state=45)
y_color = y.loc[X_color.index]
y_color = y_color.map(lambda r: 'red' if r == 1 else 'yellow')
ax = plt.axes()
ax.scatter(
    X_color.iloc[:, 0], X_color.iloc[:, 1],
    color=y_color, alpha=1)
# -----
x_axis, y_axis = np.arange(0, 1.005, .005), np.arange(0, 1.005, .005)
xx, yy = np.meshgrid(x_axis, y_axis)
xx_ravel = xx.ravel()
yy_ravel = yy.ravel()
X_grid = pd.DataFrame([xx_ravel, yy_ravel]).T
y_grid_predictions = LSVC.predict(X_grid)
y_grid_predictions = y_grid_predictions.reshape(xx.shape)
ax.contourf(xx, yy, y_grid_predictions, cmap=plt.cm.autumn_r, alpha=.3)
# -----
ax.set(
    xlabel=fields[0],
    ylabel=fields[1],
    xlim=[0, 1],
    ylim=[0, 1],
    title='decision boundary for LinearSVC');

```



### ✓ Question 3

Let's now fit a Gaussian kernel SVC and see how the decision boundary changes.

- Consolidate the code snippets in Question 2 into one function which takes in an estimator,  $x$  and  $y$ , and produces the final plot with decision boundary. The steps are:

1. fit model

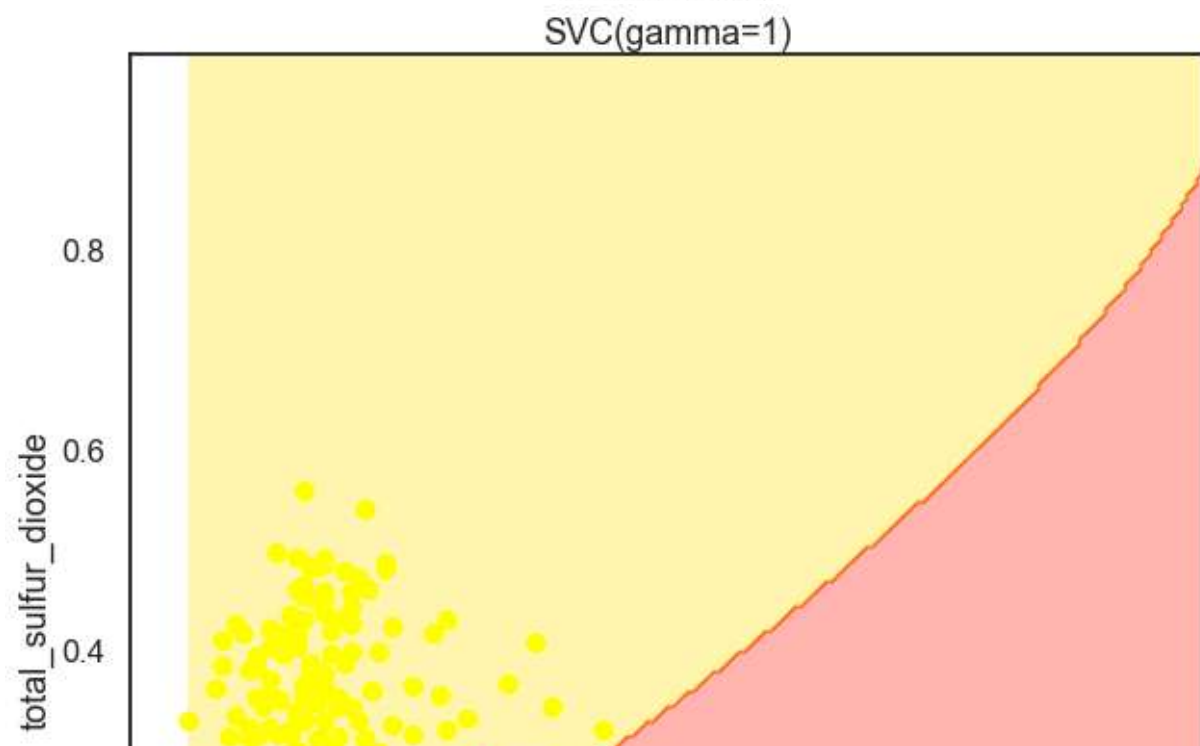
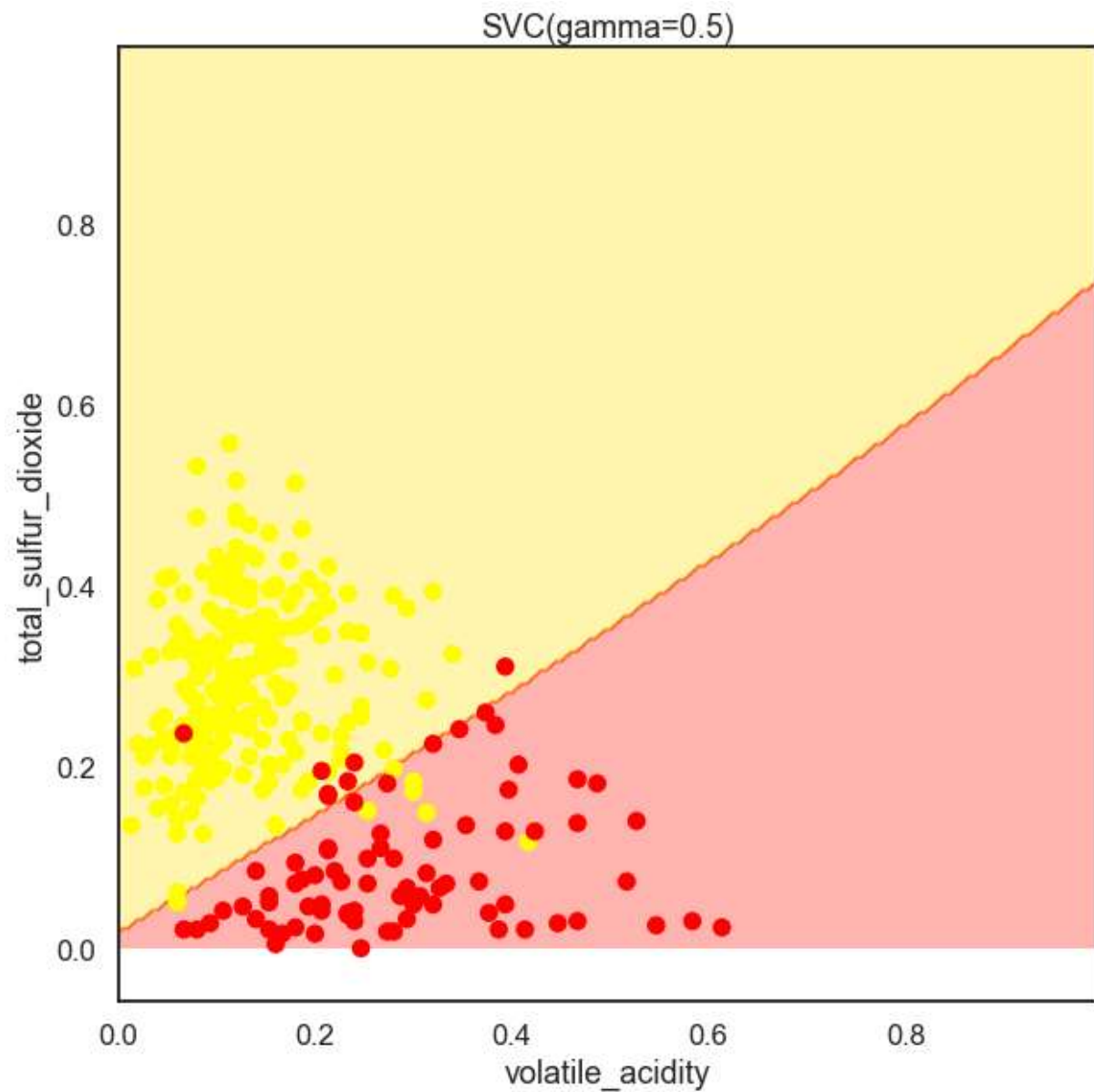
2. get sample 300 records from X and the corresponding y's
  3. create grid, predict, plot using ax.contourf
  4. add on the scatter plot
- After copying and pasting code, make sure the finished function uses your input estimator and not the LinearSVC model you built.
  - For the following values of gamma, create a Gaussian Kernel SVC and plot the decision boundary.  
gammas = [.5, 1, 2, 10]
  - Holding gamma constant, for various values of c, plot the decision boundary. You may try  
Cs = [.1, 1, 10]

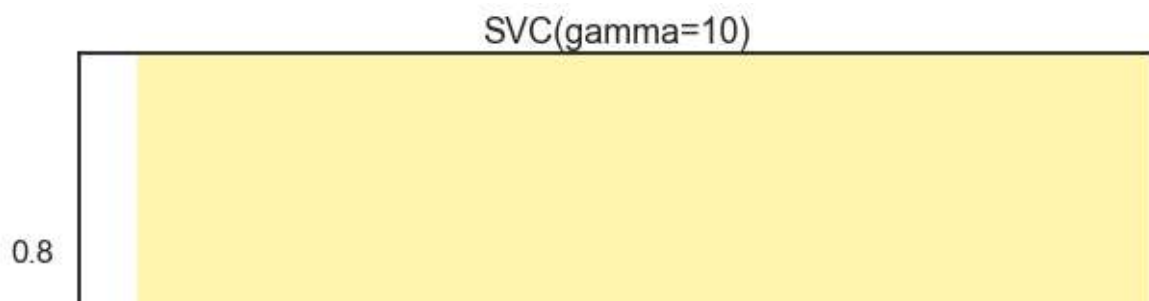
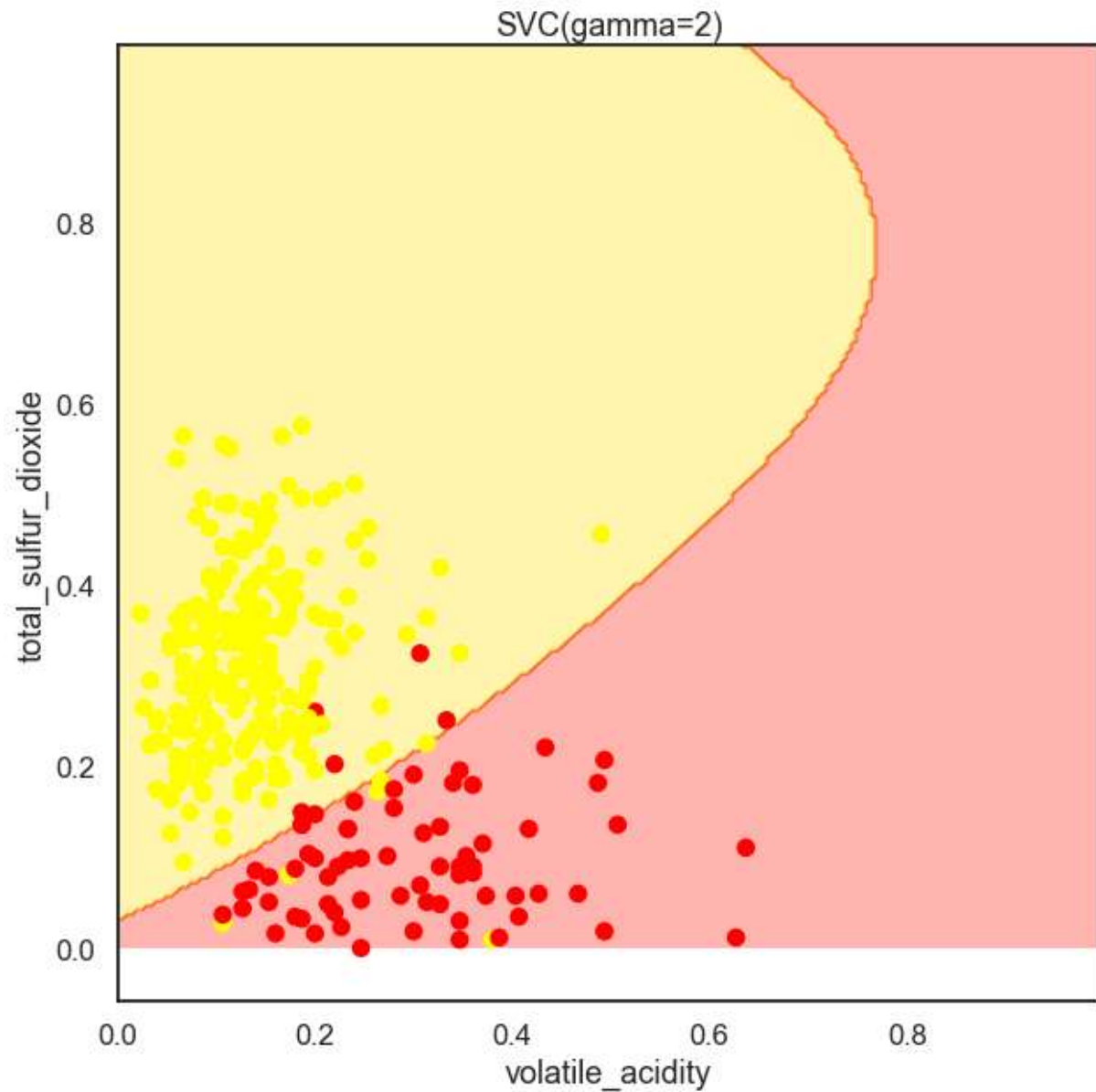
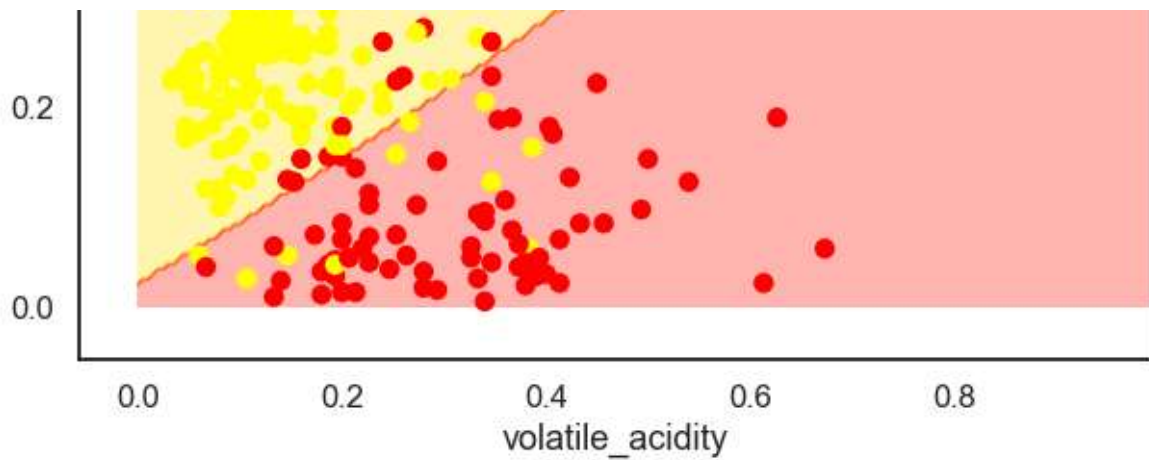
```
def plot_decision_boundary(estimator, X, y):
    estimator.fit(X, y)
    X_color = X.sample(300)
    y_color = y.loc[X_color.index]
    y_color = y_color.map(lambda r: 'red' if r == 1 else 'yellow')
    x_axis, y_axis = np.arange(0, 1, .005), np.arange(0, 1, .005)
    xx, yy = np.meshgrid(x_axis, y_axis)
    xx_ravel = xx.ravel()
    yy_ravel = yy.ravel()
    X_grid = pd.DataFrame([xx_ravel, yy_ravel]).T
    y_grid_predictions = estimator.predict(X_grid)
    y_grid_predictions = y_grid_predictions.reshape(xx.shape)

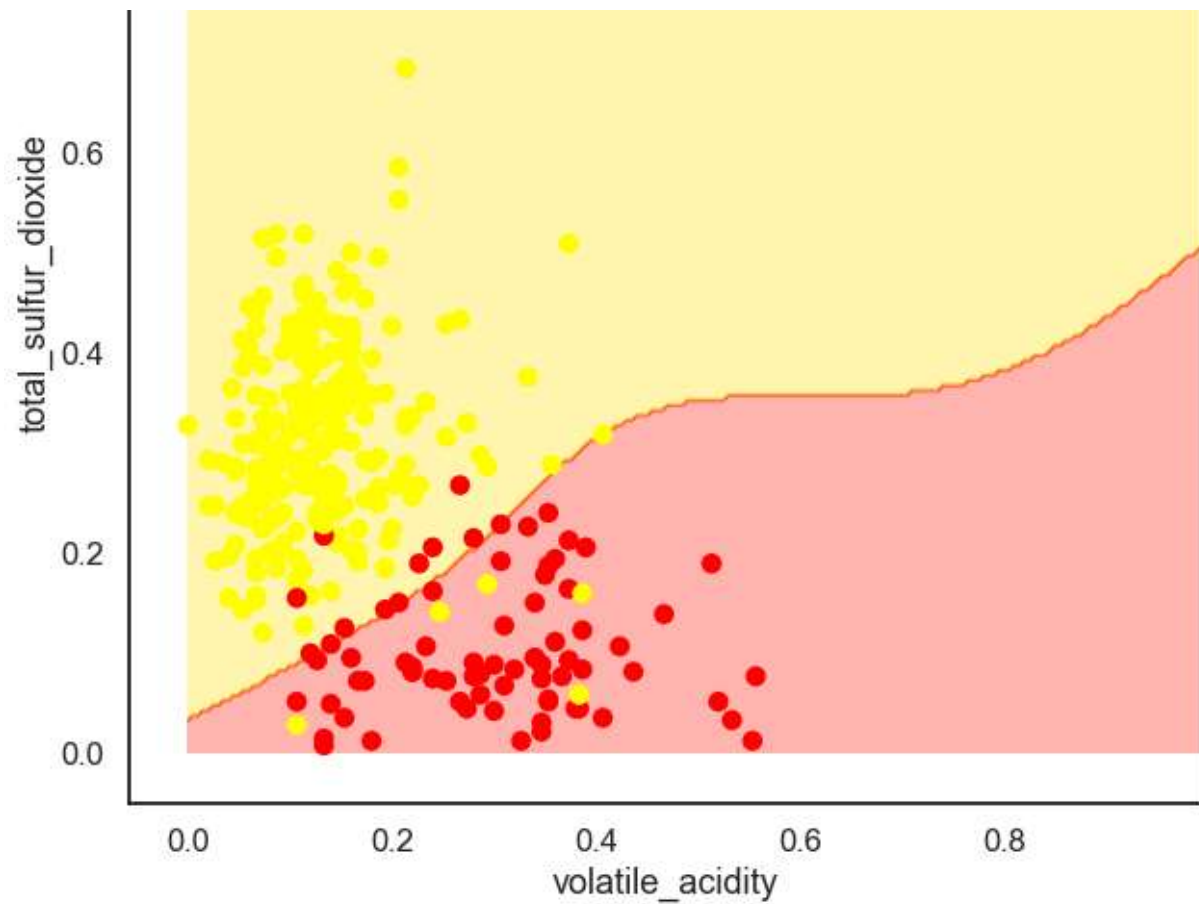
    fig, ax = plt.subplots(figsize=(10, 10))
    ax.contourf(xx, yy, y_grid_predictions, cmap=plt.cm.autumn_r, alpha=.3)
    ax.scatter(X_color.iloc[:, 0], X_color.iloc[:, 1], color=y_color, alpha=1)
    ax.set(
        xlabel=fields[0],
        ylabel=fields[1],
        title=str(estimator))

gammas = [.5, 1, 2, 10]
for gamma in gammas:
    SVC_Gaussian = SVC(kernel='rbf', gamma=gamma)
    plot_decision_boundary(SVC_Gaussian, X, y)
```









```
Cs = [.1, 1, 10]
for C in Cs:
    SVC_Gaussian = SVC(kernel='rbf', gamma=2, C=C)
    plot_decision_boundary(SVC_Gaussian, X, y)
```