fixe 0 1 2 3	A head()  ed_acidity
data fixed volat citri resid chlor free_	_sulfur_dioxide
#dat	float64 ity int64 r object e: object  fa.color equals red, and we see that that passes out true or false.  a.color == 'red'  True True True True True True True Tru
6493 6494 6495 6496 Name: # Wh #and #we' (dat	False False False False False False Color, Length: 6497, dtype: bool  Then we say astype(int) Then we see that that passes out true or false. The just going to translate each one of those trues as 1 and each one of those falses as 0.  The color == 'red').astype(int)  1 1 1 1
3 4 6492 6493 6494 6495 6496 Name: y = fiel corr	1 1  0 0 0
free_ resid citri quali alcoh pH densi fixed sulph chlor volat dtype	nol -0.032970 0.329129 ity 0.390645 d_acidity 0.486740 nates 0.487218
ax = ax.s	set_style('white')  = correlations.plot(kind='bar', color=colors[0]) set(ylim=[-1, 1], ylabel='pearson correlation');  1.0 0.5 0.0
	free_sulfur_dioxide free_sulfur_dioxide residual_sugar ditric_acid quality alcohol pH density fixed_acidity sulphates chlorides volatile_acidity
total free_ resid citri quali alcoh pH densi	want to do is get the absolute values for each.  relations.map(abs)  L_sulfur_dioxide
chlor volat dtype #We corr alcoh quali citri pH resid densi free_	rides 0.512678 tile_acidity 0.653036 e: float64  can then sort those values by calling sort values. probably just want the last two relations.map(abs).sort_values()  nol 0.032970 ity 0.119323 ic_acid 0.187397 0.329129  dual_sugar 0.34821 ity 0.348821 ity 0.390645 _sulfur_dioxide 0.471644
sulph chlor volat total dtype corr	
from  fiel prin X = scal X = prin Index Index	in sklearn.preprocessing import MinMaxScaler  ids = correlations.map(abs).sort_values().iloc[-2:].index id(fields)  data[fields] # X, rather than being all of the fields that we have available, is just going to be equal to these two fields  ler = MinMaxScaler() scaler.fit_transform(X) pd.DataFrame(X, columns=['%s_scaled' % fld for fld in fields]) #make sure that it goes back into a Pandas DataFrame rather than than a numpy array.  id(X.columns)  k(['volatile_acidity', 'total_sulfur_dioxide'], dtype='object') k(['volatile_acidity_scaled', 'total_sulfur_dioxide_scaled'], dtype='object')  Decision Boundary
<ul> <li>Fit</li> <li>Pic</li> <li>Mc</li> <li>Sc</li> <li>Us</li> <li>x_axi</li> <li>xx, y</li> <li>xx_ra</li> </ul>	that is to look at the decision boundary of a LinearSVC classifier on this dataset. Check out this example in sklearn's documentation.  It a Linear Support Vector Machine Classifier to X, y.  It is a Linear Support Ve
X_gri y_gri y_gri ax.co With Li <b>from</b> LSVC LSVC	avel = yy.ravel()  id = pd.DataFrame([xx_ravel, yy_ravel]).T  id_predictions = *[YOUR MODEL]*.predict(X_grid)  id_predictions = y_grid_predictions.reshape(xx.shape)  ontourf(xx, yy, y_grid_predictions, cmap=colors, alpha=.3)  inearSVC, it is easy to experiment with different parameter choices and see the decision boundary.  in sklearn.svm import LinearSVC  c = LinearSVC()  c.fit(X, y)  plor = X.sample(300, random_state=45)
y_co ax = ax.s	<pre>plor = y.loc[X_color.index] plor = y_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red' if r == 1 else 'yellow') plot = y.loc[X_color.map(lambda r: 'red</pre>
0.2 - 0.0 - xx.s (201,	0.0 0.2 0.4 0.6 Shape
yy.s (201, xx.r array xx.r (4040	Shape 7. 201) Favel() Favel(). shape # very big  101,)
yy.r (4040 pd.D #eve	ravel() y([0., 0., 0.,, 1., 1., 1.]) ravel().shape 01,)  pataFrame([xx_ravel, yy_ravel]).T #we put this all into a data frame, and get our grid. rry single value of y, and each of those x's and y's are meant to represent a fictional amount of either acidity or total sulfur dioxide,  0 1 0.000 0.0
1 2 3 4  40396 40397 40398	0.005 0.0 0.010 0.0 0.015 0.0 0.020 0.0
40400 40401 LSVC array #it' LSVC	rows × 2 columns  C.predict(X_grid) # single array  y([1, 1, 1,, 0, 0, 0])  S now going to be that same shape as xx as well as yy, and when we run our contour graph, we can call xx, yy.  C.predict(X_grid).reshape(xx.shape)  y([[1, 1, 1,, 1, 1],
from LSVC LSVC X_co y_co y_co y_co ax =	[1, 1, 1,, 1, 1, 1], [0, 1, 1,, 1, 1, 1],, [0, 0, 0,, 0, 0, 0], [0, 0, 0,, 0, 0, 0], [0, 0, 0,, 0, 0, 0]])  a sklearn.svm import LinearSVC  C = LinearSVC()  C.fit(x, y)  blor = X.sample(300, random_state=45)  blor = y.loc[X_color.index]  blor = y_color.map(lambda r: 'red' if r == 1 else 'yellow')  = plt.axes()
#sca # x_ax xx, xx_r yy_r X_gr y_gr y_gr ax.c	<pre>scatter( X_color.iloc[:, 0], X_color.iloc[:, 1], color=y_color, alpha=1) atter plot without the line seperating two labels  cis, y_axis = np.arange(0, 1.005, .005), np.arange(0, 1.005, .005) yy = np.meshgrid(x_axis, y_axis) ravel = xx.ravel() ravel = xy.ravel() rid = pd.DataFrame([xx_ravel, yy_ravel]).T #we put this all into a data frame, and get our grid. rid_predictions = LSVC.predict(X_grid) rid_predictions = y_grid_predictions.reshape(xx.shape) contourf(xx, yy, y_grid_predictions, cmap=plt.cm.autumn_r, alpha=.3)</pre>
	<pre>xlabel=fields[0], ylabel=fields[1], xlim=[0, 1], ylim=[0, 1], title='decision boundary for LinearSVC');</pre> <pre>.0</pre>
<b>0.</b> Gaussi	ian Kernel  ow fit a Gaussian kernel SVC and see how the decision boundary changes.
• Aff • Fo g:	onsolidate 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 for copying and pasting code, the finished function uses the input estimator and not the LinearSVC model.  or the following values of gamma, create a Gaussian Kernel SVC and plot the decision boundary.  Jammas = [.5, 1, 2, 10]  Joiding gamma constant, we plot the decision boundary for various values of C: [.1, 1, 10]  Plot_decision_boundary(estimator, X, y):
	<pre>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)</pre>
So we's initiate plot_de If you'll means reducir	ax.scatter(X_color.iloc[:, 0], X_color.iloc[:, 1], color=y_color, alpha=1) ax.set(
from gamm for	on boundary that's getting closer towards over-fitting. Perhaps, this is still a decent decision boundary. But as we keep increasing gamma, we see that we increase the complexity because we're reducing the pena ave for our regularization terms.  In sklearn.svm import SVC  In sklearn.svm imp
	.8-
total	.2-
	0.0 0.2 0.4 0.6 0.8 volatile_acidity  SVC(gamma=1)
total_sulfur_dioxide	.6 -
	.0.0 0.2 0.4 0.6 0.8 volatile_acidity
	SVC(gamma=2) .8-
total_sulfur_dioxide	.4-
0.	0.0 0.2 0.4 0.6 0.8 volatile_acidity  SVC(gamma=10)
_sulfur_dioxide	.6-
o. total	
we're g being to remem Then th regular  Cs = for	volatile_acidity  going to run through each one of our Cs in just a second. So with the Cs, we're going to have, again, increasing from 0.1, to 1, to 10, and we're going to call SVC. We're going to say gamma equal to a constant, he two, then we're going to loop through each one of these different C values. Again, for that I want you to start thinking which ones are going to have higher versus lower values in regards to the amount of regularization ber, more regularization for, I want you to think for higher value of C versus lower value of C, as well as think about the complexity of the model once you reduce or increase the regularization term.  the same holds for C. So it's the same inverse relationship where the lower value of C, so here we're starting off with C equals 0.1, lower value of C means that we're going to have higher regularization. Higher rization means a less complex model. Then as we increase C, we see how it gets more complex, more curvature to our model.  [1, 1, 10]  C in Cs:  SVC_Gaussian = SVC(kernel='rbf', gamma=2, C=C)  plot_decision_boundary(SVC_Gaussian, X, y)
SV	plot_decision_boundary(SVC_Gaussian, X, y)  VC(C=0.1, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=3, gamma=2, kernel='rbf', max_iter=-1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False)
total_sulfur_dioxide	.6
0.	2 0.0 0.2 0.4 0.6 0.8 volatile_acidity
0.	SVC(C=1, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=3, gamma=2, kernel='rbf', max_iter=-1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False)
total_sulfur_dioxide	
0. 0. S\	.0 0.0 0.2 0.4 0.6 0.8 volatile_acidity
0.	VC(C=10, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=3, gamma=2, kernel='rbf', max_iter=-1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False)
total_sulfur_dioxide	.4
	0.0 0.2 0.4 0.6 0.8 1.0 volatile_acidity  aring Kernel Execution Times
Jupyter We pro  Cr Us	exercise, we will compare the fitting times between SVC vs Nystroem with rbf kernel.  In Notebooks provide a useful magic function <code>%timeit</code> which executes a line and prints out the time it took to fit. If we type <code>%%timeit</code> in the beginning of the cell, it will output the execution time. Deceded with the following steps:  In reate <code>y</code> from data.color, and <code>X</code> from the rest of the columns.  See <code>%%timeit</code> to get the time for fitting an SVC with rbf kernel.  See <code>%%timeit</code> to get the time for the following: fit_transform the data with Nystroem and then fit a SGDClassifier.  See <code>%%timeit</code> to get the time for the following: fit_transform the data with Nystroem and then fit a SGDClassifier.
<ul> <li>Ma</li> <li>Ma</li> <li>Co</li> <li>from from</li> <li>from</li> <li>x =</li> <li>kwar</li> <li>svc</li> </ul>	ake 5 copies of X and concatenate them ake 5 copies of y and concatenate them compare the time it takes to fit the both methods above  in sklearn.kernel_approximation import Nystroem in sklearn.svm import SVC in sklearn.linear_model import SGDClassifier  data.color == 'red' data[data.columns[:-1]]  rgs = {'kernel': 'rbf'} #kwargs, and we can parse that in as just a something that you pass him within Python. As long as you pass in a dictionary, you're saying for the = SVC(**kwargs)
svc nyst sgd %%ti svc. 632 m	= SVC(**kwargs) rroem = Nystroem(**kwargs) = SGDClassifier()  Immeit fit(X, y)  Ins ± 40.5 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)  Fixed_acidity volatile_acidity citric_acid residual_sugar chlorides free_sulfur_dioxide total_sulfur_dioxide density pH sulphates alcohol quality  7.4 0.70 0.00 1.9 0.076 11.0 34.0 0.99780 3.51 0.56 9.4 5
1 2 3 4  6492 6493 6494 6495	
	6.0 0.21 0.38 0.8 0.020 22.0 98.0 0.98941 3.26 0.32 11.8 6  ows × 12 columns