Figure 1: Plot of different SVC classification

# 1 Chosen kernel: RBF(C=2, Gamma=0.1)

1. Used ScikitLearn SVM package for this classification process. Module has 4 different kernel, namely : Linear, Polynomial, RBF , Sigmoid. These four kernels have multiple hyperparameters which can be tuned to get the optimal classifier which avoids overfitting.

2. Understanding the influence of hyperparamter 'C' is crucial for understanding if our model has generalizability or not. Following is the SVM primal equations:

$$
\begin{aligned}
\min_{w,b,\zeta} \frac{1}{2} w^T w + C \sum_{i=1}^{n} \zeta_i \\
\text{subject to } y_i(w^T \phi(x_i) + b) \geq 1 - \zeta_i, \\
\zeta_i \geq 0, i = 1, ..., n
\end{aligned}
\tag{1}
$$

C acts as a regularization parameter for the soft error.C not too large: the outliers wont affect the decision boundary much. C large: SVM will try to (over)fit the decision boundary to the outliers.

3. GridSearch CV was used to iterate over different possible combinations of Kernels with their hyperparameters and 5 fold cross-validation was done and accuracy measure was used to pick the best models.

4. Figure below shows the best accuracy measure giving kernels. Linear Kernel accuracy: 0.823 (+-0.058) , RBF C=2 accuracy: 0.864 (+-0.07), RBF C=10 accuracy: 0.88 (+-0.06). From these we can see that RBF kernel with C=10 has the best accuracy but has large C value which may cause overfitting and have less generalization ability. So I used the C=2 and gamma=0.1 RBF kernel which has a simpler model and similar accuracy values.