

Supervised Learning - Support Vector Machines - Lab 06

In []:

Name = "Muhammad Omer Farooq Bhatti"

Id = "st122498"

Mathematical Details

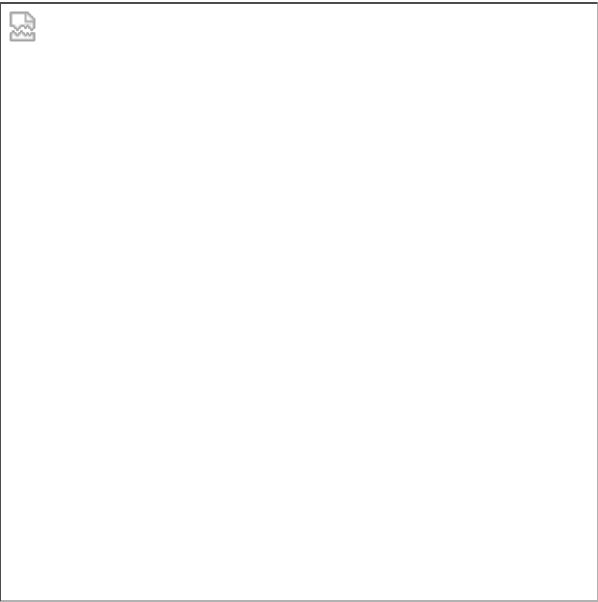
Notations

Scalars are denoted with italic lowercases (e.g., y , b), vectors with bold lowercases (e.g., \mathbf{w} , \mathbf{x}), and matrices with italic uppercases (e.g., W). \mathbf{w}^T is the transpose of \mathbf{w} and $\|\mathbf{w}\| = \sqrt{\mathbf{w}^T \mathbf{w}}$

Let:

- \mathbf{x} be a feature vector (i.e., the input of the SVM). $\mathbf{x} \in \mathbb{R}^n$, where n is the dimension of the feature vector.
- y be the class (i.e., the output of the SVM). $y \in \{-1, 1\}$, i.e. the classification task is binary.
- \mathbf{w} and b be the parameters of the SVM: we need to learn them using the training set.
- $(\mathbf{x}^{(i)}, y^{(i)})$ be the i th sample in the dataset. Let's assume we have m samples in the training set.

With $n = 2$, one can represent the SVM's decision boundaries as follows:



The hypothesis function is simply:

$$h = \text{sign}(\mathbf{w}^T \mathbf{x}^{(i)} + b)$$

The constraint is as follows:

$$y^{(i)} = \begin{cases} -1 & \text{if } \mathbf{w}^T \mathbf{x}^{(i)} + b \leq -1 \\ 1 & \text{if } \mathbf{w}^T \mathbf{x}^{(i)} + b \geq 1 \end{cases}$$

If we multiply y on both equations, the constraint function can then be more concisely written as

$$y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1$$

Goal

The SVM aims at satisfying two requirements:

- The SVM should maximize the distance between the two decision boundaries. Mathematically, this means we want to maximize the distance between the hyperplane defined by $\mathbf{w}^T \mathbf{x} + b = -1$ and the hyperplane defined by $\mathbf{w}^T \mathbf{x} + b = 1$. This distance is equal to $\frac{2}{\|\mathbf{w}\|}$. This means we want to solve $\max_{\mathbf{w}} \frac{2}{\|\mathbf{w}\|}$. Equivalently we want $\min_{\mathbf{w}} \frac{\|\mathbf{w}\|}{2}$
- The SVM should also correctly classify all $\mathbf{x}^{(i)}$, which means

$$y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1, \forall i \in \{1, \dots, m\}$$

Primal Problem

Recall that our primal, optimization problem is of the form:

$$\begin{aligned} \min_{w,b} f(w,b) &= \min_{w,b} \frac{1}{2} \|w\|^2 \\ s.t. \quad g_i(w,b) &= -(y^{(i)}(w^T x^{(i)} + b) - 1) \leq 0 \end{aligned}$$

Lagrange Method

The method of Lagrange multipliers allows us to turn a constrained optimization problem into an unconstrained one of the form.

$$\mathcal{L} = f(w) + \sum \alpha \cdot g(w) + \sum \beta \cdot h(w)$$

where $f(w)$ is the function we want to minimize, and $g_i(w)$ and $h_i(w)$ is the constraint function where $g_i(w) \leq 0$ and $h_i(w) = 0$. To

solve this function, we set the partial derivative of our parameters to 0.

We can write our equation as

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_i^m \alpha^{(i)} [y^{(i)} (w^T x^{(i)} + b) - 1]$$

Where $\mathcal{L}(w, b, \alpha)$ is called the Lagrangian and α_i are called the Lagrangian multipliers.

Our primal optimization problem with the Lagrangian becomes the following:

$$\min_{w, b} \left(\max_{\alpha} \mathcal{L}(w, b, \alpha) \right)$$

Dual Problem

This is the idea of turning primal problem into dual problem by acknowledging that this is roughly the same:

$$\min_{w, b} \left(\max_{\alpha} \mathcal{L}(w, b, \alpha) \right) = \max_{\alpha} \left(\min_{w, b} \mathcal{L}(w, b, \alpha) \right)$$

Why Dual, not Primal?

Short answer: Faster computation + allows to use the kernel trick.

Duality and KKT

Karush Kuhn Tucker (KKT) conditions allow us to solve the dual problem instead of the primal one, while ensuring that the optimal solution is the same. In our case the conditions are the following:

- The primal objective and inequality constraint functions must be convex
- The equality constraint function must be affine
- The constraints must be strictly feasible

Then there exists w^*, α^* which are solutions to the primal and dual problems. Moreover, the parameters w^*, α^* satisfy the KKT conditions below:

$$\frac{\partial}{\partial w_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0 \quad (A)$$

$$\frac{\partial}{\partial \beta_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0 \quad (B)$$

$$\alpha_i^* g_i(w^*) = 0 \quad (C)$$

$$g_i(w^*) \leq 0 \quad (D)$$

$$\alpha_i^* \geq 0 \quad (E)$$

Moreover, if some w^*, α^* satisfy the KKT solutions then they are also solution to the primal and dual problem.

Note that when $g_i(w^*) < 0$ (samples that are not in the gutter), to maximize, α must be 0, in order to satisfy C . (Note that our equation uses minus sign instead of plus sign, but the result is the same.). This means that any samples having $\alpha = 0$ is not support vectors. On the other hand, when $g_i(w^*) = 0$ (samples that are in the gutter), to maximize, α can or cannot be 0 to satisfy C , that is, we can choose whether to choose that sample to be our support vector. Any samples with $\alpha > 0$ is support vector.

Back to the Dual Problem

Now that you are convinced dual problem is worth solving more than primal problem. To solve dual problem, we take the partial derivatives of $\mathcal{L}(w, b, \alpha)$ with respect to w and b , equate to zero (according to KKT conditions) and then plug the results back into the original equation of the Lagrangian

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_i^m \alpha_i [y^{(i)} (w^T x^{(i)} + b) - 1]$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{w} - \sum \alpha_i y^{(i)} \mathbf{x}^{(i)} = 0$$

$$\mathbf{w} = \sum \alpha_i y^{(i)} \mathbf{x}^{(i)} \quad (6)$$

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum \alpha_i y^{(i)} = 0,$$

$$\sum \alpha_i y^{(i)} = 0 \quad (7)$$

Plugging the w back to the original equation, as well as the fact that $\sum \alpha_i y^{(i)} = 0$, we got that

$$\mathcal{L}(w, b, \alpha) = \sum_i^m \alpha_i - \frac{1}{2} \sum_{i,j}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)} x^{(j)} \rangle$$

hence generating an equivalent dual optimization problem of the form

$$\begin{aligned} & \max_{\alpha} \min_{w,b} \mathcal{L}(w, b, \alpha) \\ & \max_{\alpha} \sum_i^m \alpha_i - \frac{1}{2} \sum_{i,j}^m y^{(i)} y^{(j)} \alpha_i \alpha_j < x^{(i)} x^{(j)} > \\ & s. t. \quad \alpha_i \geq 0 \\ & s. t. \quad \sum_i^m \alpha_i y^{(i)} = 0 \end{aligned}$$

What's next

To max α according to the equation, we cannot use gradient since the equation is neither concave nor convex. But notice the quadratic form of the equation. To optimize such equation, we can use **quadratic programming**. Later in the course, we shall discuss using **CVXOPT** python library. Solving the equation gives us back α , which using the α , we can get w , and using w , we can then get b . Hooray, we got our hypothesis function $\mathbf{w}^T \mathbf{x}^{(i)} + b$

We shall also discuss other minor improvements to the SVM. Namely two - 1) **Soft margin**, 2) **Kernels**.

Hard margin vs. Soft margin

What we have done so far is the **hard-margin SVM**, as this quadratic optimization problem admits a solution iff the data is linearly separable.

One can relax the constraints by introducing so-called **slack variables** $\xi^{(i)}$. Note that each sample of the training set has its own slack variable. This gives us the **soft-margin SVM** following quadratic optimization problem:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{\|\mathbf{w}\|^2}{2} + C \sum_{i=1}^m \xi^{(i)}, \\ s. t. \quad & y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \xi^{(i)}, \quad \forall i \in \{1, \dots, m\} \\ & \xi^{(i)} \geq 0, \quad \forall i \in \{1, \dots, m\} \end{aligned}$$

Following the inequality constraints on the slack variables, we add a beta in front of it (following the Lagrange method). Thus the Lagrange formulation change to:

$$\mathcal{L}(w, b, \xi, \alpha, \beta) = \frac{1}{2} ||w||^2 + C \sum_{i=1}^m \xi^{(i)} - \sum_i^m \alpha_i [y^{(i)} (w^T x^{(i)} + b) - 1 - \xi^{(i)}] - \sum_{i=1}^m \beta_i \xi^{(i)}$$

Although this looks tough, terms cancel out nicely.

Following dual form, we need to first minimize w.r.t \mathbf{w} , b , and ξ , before maximizing w.r.t each $\alpha_i \geq 0$ and $\beta_i \geq 0$:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{w}} &= \mathbf{w} - \sum \alpha_i y^{(i)} \mathbf{x}_i = 0 \\ \mathbf{w} &= \sum \alpha_i y^{(i)} \mathbf{x}^{(i)} \\ \frac{\partial \mathcal{L}}{\partial b} &= - \sum \alpha_i y^{(i)} = 0, \\ \sum \alpha_i y^{(i)} &= 0 \\ \frac{\partial \mathcal{L}}{\partial \xi} &= C - \sum \alpha_i - \sum \beta_i = 0 \end{aligned}$$

Plugging these truths to the Lagrange formula, it happens that we got the same formula, with one addition, i.e., $\alpha_i \leq C$, thus

$$\begin{aligned} & \max_{\alpha} \min_{w,b,\xi} \mathcal{L}(w, b, \xi, \alpha, \beta) \\ & \max_{\alpha} \sum_i^m \alpha_i - \frac{1}{2} \sum_{i,j}^m y^{(i)} y^{(j)} \alpha_i \alpha_j < x^{(i)} x^{(j)} > \\ & s. t. \quad C \geq \alpha_i \geq 0 \\ & s. t. \quad \sum_i^m \alpha_i y^{(i)} = 0 \end{aligned}$$

Kernels

Recall this equation:

$$\mathbf{w} = \sum \alpha_i y^{(i)} \mathbf{x}^{(i)}$$

Attempting to plug this back to our hypothesis function, we get

$$h = y^{(i)} (\sum \alpha_i y^{(i)} \mathbf{x}^{(i)} \mathbf{x}'^{(i)} + b) \geq 1$$

Here we put $'$ to another x to denote that these two x are different. Make an observation that to predict anything, we simply need to find the dot product of a pair of x

Based on this observation, one can add even more flexibility by introducing a function ϕ that maps the original feature space to a

higher dimensional feature space. This allows non-linear decision boundaries.

Suppose we have a mapping $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^d$, where d is of higher dimensions. Then the dot product of \mathbf{x} and \mathbf{x}' in this space is $\varphi(\mathbf{x})^T \varphi(\mathbf{x}')$. A kernel is a function k that corresponds to this dot product without having to performing φ , i.e.,

$$k(\mathbf{x}, \mathbf{x}') = \varphi(\mathbf{x})^T \varphi(\mathbf{x}')$$

Why is this useful? or How is this different from previous transformations such as polynomial feature engineering we did on linear regression? That is, kernels give a way to **compute dot products in some feature space without even knowing what this space is and what is φ** . Since we don't need to actually transform features, this saves a lot of computation time.

For example, consider a simple polynomial kernel of degree 2, $k(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^T \mathbf{x}')^2$ (Note that in general polynomial kernel, 2 can be p) with $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^2$. This doesn't seem to correspond to any mapping function φ . It's just a function that returns a real number. Assuming that $\mathbf{x} = (x_1, x_2)$ and $\mathbf{x}' = (x'_1, x'_2)$, let's expand this expression:

$$\begin{aligned} k(\mathbf{x}, \mathbf{x}') &= (1 + \mathbf{x}^T \mathbf{x}')^2 = (1 + x_1 x'_1 + x_2 x'_2)^2 = \\ &= 1 + x_1^2 x'^2_1 + x_2^2 x'^2_2 + 2x_1 x'_1 + 2x_2 x'_2 + 2x_1 x_2 x'_1 x'_2 \end{aligned}$$

Note that this is nothing else but a dot product between two vectors of $\varphi(\mathbf{x}) = (1, x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2)$ and $\varphi(\mathbf{x}') = (1, x'^2_1, x'^2_2, \sqrt{2}x'_1, \sqrt{2}x'_2, \sqrt{2}x'_1x'_2)$. So the kernel $k(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^T \mathbf{x}')^2 = \varphi(\mathbf{x})^T \varphi(\mathbf{x}')$ computes a dot product in 6-dimensional space without explicitly visiting this space.

Another example is Gaussian kernel:

$$k(\mathbf{x}, \mathbf{x}') = \exp \left(- \gamma \| \mathbf{x} - \mathbf{x}' \|^2 \right)$$

Here is a short proof, ignoring γ , we can use the taylor series expansion, then we get:

$$\begin{aligned} &= \exp \left(- \| \mathbf{x} - \mathbf{x}' \|^2 \right) \\ &= \exp \left(- \mathbf{x}^2 \right) \exp \left(- \mathbf{x}'^2 \right) \sum_{k=0}^{\infty} \frac{2^k (\mathbf{x})^k (\mathbf{x}')^k}{k!} \end{aligned}$$

Using kernel, our optimization problem can be revised to:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{\|\mathbf{w}\|}{2} + C \sum_{i=1}^m \xi^{(i)}, \\ s. \ t. \quad & y^{(i)}(\mathbf{w}^T \phi(\mathbf{x}^{(i)}) + b) \geq 1 - \xi^{(i)}, \quad \forall i \in \{1, \dots, m\} \\ & \xi^{(i)} \geq 0, \quad \forall i \in \{1, \dots, m\} \end{aligned}$$

And our lagrange formula changed only slightly to

$$\begin{aligned} &\max_{\alpha} \min_{w, b, \xi} \mathcal{L}(w, b, \xi, \alpha, \beta) \\ &\max_{\alpha} \sum_i^m \alpha_i - \frac{1}{2} \sum_{i,j}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \varphi(\mathbf{x}^{(i)}) \varphi(\mathbf{x}^{(j)}) \\ s. \ t. \quad & C \geq \alpha_i \geq 0 \\ s. \ t. \quad & \sum_i^m \alpha_i y^{(i)} = 0 \end{aligned}$$

Making a prediction

Once the α_i are learned, one can predict the class of a new sample with the feature vector `X_test` as follows:

1. Scratch: SVM with Hard Margin

Since we will solve this optimization problem using the **CVXOPT** library in python, we will need to match the solver's API which, according to the documentation is of the form:

$$\begin{aligned} \min \quad & \frac{1}{2} x^T P x + q^T x \\ s. \ t. \quad & G x \leq h \\ & A x = b \end{aligned}$$

Recall that the dual problem is expressed as:

$$\max_{\alpha} \sum_i^m \alpha_i - \frac{1}{2} \sum_{i,j}^m y^{(i)} y^{(j)} \alpha_i \alpha_j < x^{(i)} x^{(j)} >$$

Let **H** be a matrix such that

$$H_{ij} = y^{(i)} y^{(j)} < x^{(i)} x^{(j)} >$$

Then the optimization becomes:

$$\begin{aligned} \max_{\alpha} \quad & \sum_i^m \alpha_i - \frac{1}{2} \alpha^T \mathbf{H} \alpha \\ \text{s.t.} \quad & \alpha_i \geq 0 \\ & \sum_i^m \alpha_i y^{(i)} = 0 \end{aligned}$$

We convert the sums into vector form and multiply both the objective and the constraint by -1 which turns this into a minimization problem and reverses the inequality

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T \mathbf{H} \alpha - 1^T \alpha \\ \text{s.t.} \quad & -\alpha_i \leq 0 \\ & y^T \alpha = 0 \end{aligned}$$

We have that strange $1.T$ just for the sake of matching with `CVXOPT` solver API.

We are now ready to convert our numpy arrays into the `cvxopt` format, using the same notations as in the documentation.

- $\mathbf{P} = \mathbf{H}$ a matrix of size (m, m)
- $\mathbf{q} = -\mathbf{1}$ a matrix of size $(m, 1)$
- $\mathbf{G} = -\text{diag}[\mathbf{1}]$ a diagonal matrix of -1 s size (m, m)
- $\mathbf{h} = \mathbf{0}$ a vector of size $(m, 1)$
- $\mathbf{A} = \mathbf{y}$ a vector of size $(m, 1)$
- $\mathbf{b} = \mathbf{0}$ a scalar

Note that in the simple example of $m = 2$ the matrix \mathbf{G} and vector \mathbf{h} which define the constraint are

$$G = \begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix}$$
$$h = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Computing the matrix H in vectorized form

Consider the simple example with 2 input samples $\{x^{(1)}, x^{(2)}\} \in \mathbb{R}^2$ which are two dimensional vectors

$$X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} \\ x_1^{(2)} & x_2^{(2)} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \end{bmatrix}$$

We now proceed to create a new matrix X' where each input sample x is multiplied by the corresponding output label y . This can be done easily in Numpy using vectorization and padding.

$$X' = \begin{bmatrix} x_1^{(1)} y^{(1)} & x_2^{(1)} y^{(1)} \\ x_1^{(2)} y^{(2)} & x_2^{(2)} y^{(2)} \end{bmatrix}$$

Finally we take the **matrix multiplication** of X' and its transpose giving $H = X' X'^T$

2. Scratch: SVM with Soft Margin

To make the data no longer linearly separable, we shall add a positive point in the middle of the negative cluster:

For the softmax margin SVM, the optimization problem can be slightly revised and be expressed as

$$\begin{aligned} \max_{\alpha} \quad & \sum_i^m \alpha_i - \frac{1}{2} \alpha^T \mathbf{H} \alpha \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C \\ & \sum_i^m \alpha_i y^{(i)} = 0 \end{aligned}$$

which can be written in standard form as

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T \mathbf{H} \alpha - 1^T \alpha \\ \text{s.t.} \quad & -\alpha_i \leq 0 \\ & \alpha_i \leq C \\ & y^T \alpha = 0 \end{aligned}$$

We translate this new constraint into standard form by concatenating below matrix \mathbf{G} a diagonal matrix of 1s of size (m, m) . Similarly for the vector \mathbf{h} to which the value of \mathbf{C} is added m times.

Note that in the simple example of $m = 2$ the matrix \mathbf{G} and vector h which define the constraint are:

$$G = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$$
$$\begin{bmatrix} 0 \end{bmatrix}$$

When to Use SVM

SVM work fantastic, given the following advantages:

- Kernel trick allows non-linearity and works well with high-dimensional data
- Supports regularization out of the box with the slack variables
- Use only few support vectors as decision boundary thus the prediction phase is very fast.

However, SVMs have several disadvantages as well:

- For large number of training amples, the computational cost can be scary! Around $\mathcal{O}[N^3]$ at worst
- Require fine-tuning C which often takes a lot of time!
- Does not support probability - only output the class

We usually use SVM when our Baye or Logistic fellows fail us, or we believe our data is non-linear to some extent.

===Task===

Your work:

- Load this dataset to numpy, with first two columns as features and last as target
- Plot the data using a scatter plot
- Perform the SVM classification using our scratch code; put it into class, and allows users to use soft or hard margin

```
In [32]: import numpy as np
import matplotlib.pyplot as plt
import cvxopt
from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report
```

In [39]:

```

class SVM:
    def __init__(self, kernel='gaussian', C=1):
        self.kernel = kernel
        self.C = C

    def fit(self, X, y, margin='hard'):
        self.y = y
        self.X = X
        m, n = X.shape
        self.margin=margin

        # Calculate Kernel
        self.K = np.zeros((m, m))
        for i in range(m):
            self.K[i, :] = self.calc_similarity_functions(X[i, np.newaxis], self.X)

        # Solve with cvxopt final QP needs to be reformulated
        # to match the input form for cvxopt.solvers.qp
        P = cvxopt.matrix(np.outer(y, y) * self.K)
        q = cvxopt.matrix(-np.ones((m, 1)))
        if self.margin=='hard':
            G = cvxopt.matrix(np.eye(m) * -1)
            h = cvxopt.matrix(np.zeros(m))
        elif self.margin=='soft':
            G = cvxopt.matrix(np.vstack((np.eye(m) * -1, np.eye(m))))
            h = cvxopt.matrix(np.hstack((np.zeros(m), np.ones(m) * self.C)))
        else:
            raise ValueError("Invalid margin value given.")

        A = cvxopt.matrix(y, (1, m), "d")
        b = cvxopt.matrix(np.zeros(1))
        cvxopt.solvers.options["show_progress"] = False
        sol = cvxopt.solvers.qp(P, q, G, h, A, b)
        self.alphas = np.array(sol["x"])

    def calc_similarity_functions(self, x, z, sigma=0.1, p=5):
        if self.kernel=='gaussian':
            return np.exp(-np.linalg.norm(x - z, axis=1) ** 2 / (2 * (sigma ** 2)))
        elif self.kernel=='polynomial':
            return (1 + np.dot(x, z.T)) ** p
        elif self.kernel=='linear':
            return np.dot(x, z.T)
        else:
            raise ValueError("Invalid entry for kernel")

    def predict(self, X): #<----this is X_test
        y_predict = np.zeros((X.shape[0]))
        sv = self.get_parameters(self.alphas)

        for i in range(X.shape[0]):
            y_predict[i] = np.sum(
                self.alphas[sv]
                * self.y[sv, np.newaxis]
                * self.calc_similarity_functions(X[i], self.X[sv])[:, np.newaxis]
            )

        return np.sign(y_predict + self.b)

    def get_parameters(self, alphas):
        threshold = 1e-5
        if self.margin=='hard':
            sv = (alphas > threshold).flatten()
        elif self.margin=='soft':
            sv = ((alphas > threshold) * (alphas < self.C)).flatten()
        self.w = np.dot(self.X[sv].T, alphas[sv] * self.y[sv, np.newaxis])
        self.b = np.mean(
            self.y[sv, np.newaxis]
            - self.alphas[sv] * self.y[sv, np.newaxis] * self.K[sv, sv][:, np.newaxis]
        )
        return sv

    def plot_contour(self):
        # plot the resulting classifier
        h = 0.01
        x_min, x_max = self.X[:, 0].min() - 1, self.X[:, 0].max() + 1
        y_min, y_max = self.X[:, 1].min() - 1, self.X[:, 1].max() + 1

        xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))

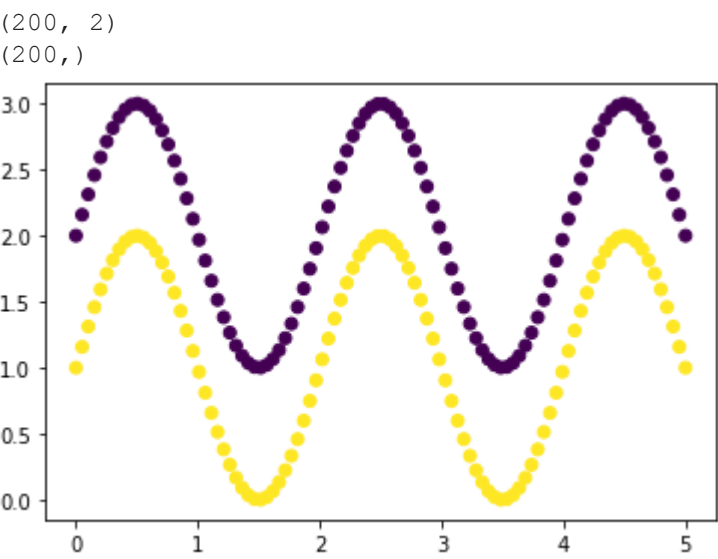
        points = np.c_[xx.ravel(), yy.ravel()]

        Z = self.predict(points)
        Z = Z.reshape(xx.shape)
        plt.contourf(xx, yy, Z, cmap=plt.cm.Spectral, alpha=0.8)

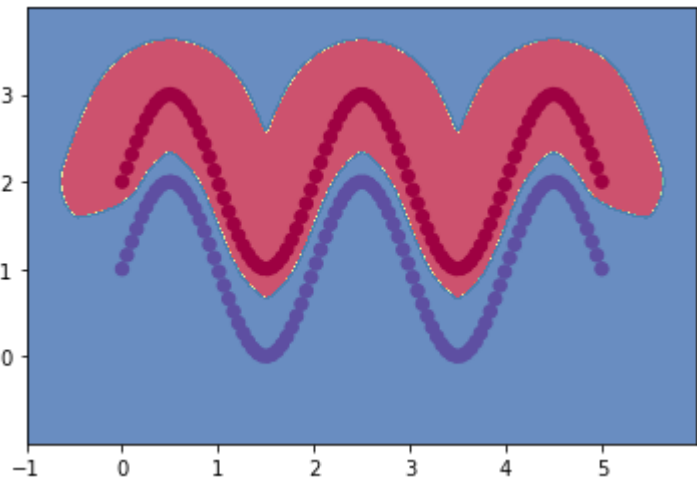
        # plt the points
        plt.scatter(X[:, 0], X[:, 1], c=y, s=40, cmap=plt.cm.Spectral)

```

```
In [30]: dataset = [[3.63636364e+00,1.09036800e+00,0], [4.09090909e+00,2.28173256e+00,0], [1.01010101e-01,1.31203345e
dataset = np.array(dataset)
X = dataset[:, :2]
y = dataset[:, -1]
y[y==0] = -1
print(X.shape)
print(y.shape)
plt.scatter(X[:,0], X[:,1], c=y, label="Features of our data")
plt.show()
X_train, X_test, y_train, y_test = train_test_split(X,y, test_size = 0.4)
```



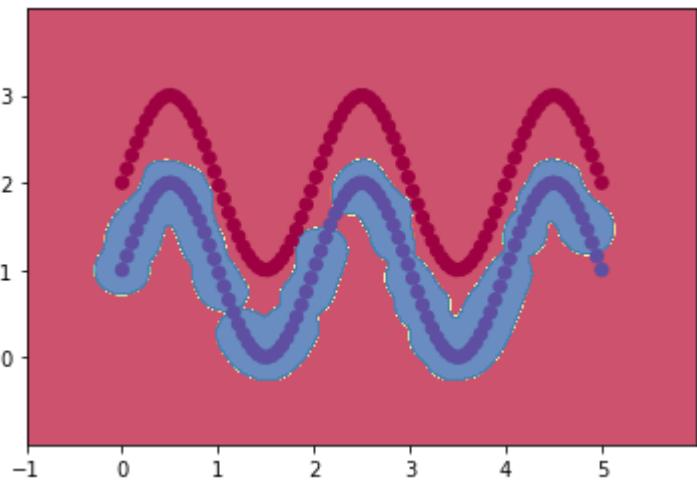
```
In [51]: svm = SVM(kernel='gaussian',C=1)
svm.fit(X, y, margin='hard')
svm.plot_contour()
```



```
In [40]: svm = SVM(kernel='gaussian',C=1)
svm.fit(X_train, y_train, margin='hard')
y_pred = svm.predict(X_test)
print(classification_report(y_test, y_pred))
svm.plot_contour()

#print(f"Accuracy: {sum(y==y_pred)/y.shape[0]}")
```

	precision	recall	f1-score	support
-1.0	0.91	1.00	0.95	39
1.0	1.00	0.90	0.95	41
accuracy			0.95	80
macro avg	0.95	0.95	0.95	80
weighted avg	0.95	0.95	0.95	80




```
In [50]: for kernel in ['linear', 'polynomial', 'gaussian']:
        for c in [1,10]:
            svm = SVM(kernel=kernel,C=c)
            svm.fit(X_train, y_train, margin='soft')
            y_pred = svm.predict(X_test)
            print("For kernel = ", kernel)
            print("For C = ", c)
            print(classification_report(y_test, y_pred))

For kernel = linear
For C = 1
      precision    recall  f1-score   support

   -1.0         0.49      1.00      0.66         39
    1.0         0.00      0.00      0.00         41

 accuracy          0.24
macro avg          0.24      0.50      0.33         80
weighted avg          0.24      0.49      0.32         80

For kernel = linear
For C = 10
      precision    recall  f1-score   support

   -1.0         0.49      1.00      0.66         39
    1.0         0.00      0.00      0.00         41

 accuracy          0.24
macro avg          0.24      0.50      0.33         80
weighted avg          0.24      0.49      0.32         80

For kernel = polynomial
For C = 1
      precision    recall  f1-score   support

   -1.0         0.00      0.00      0.00         39
    1.0         0.51      1.00      0.68         41

 accuracy          0.26
macro avg          0.26      0.50      0.34         80
weighted avg          0.26      0.51      0.35         80

For kernel = polynomial
For C = 10
      precision    recall  f1-score   support

   -1.0         0.00      0.00      0.00         39
    1.0         0.51      1.00      0.68         41

 accuracy          0.26
macro avg          0.26      0.50      0.34         80
weighted avg          0.26      0.51      0.35         80

For kernel = gaussian
For C = 1
C:\pythonDSAI\lib\site-packages\sklearn\metrics\_classification.py:1248: UndefinedMetricWarning: Precision a
nd F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` par
ameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
C:\pythonDSAI\lib\site-packages\sklearn\metrics\_classification.py:1248: UndefinedMetricWarning: Precision a
nd F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` par
ameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
C:\pythonDSAI\lib\site-packages\sklearn\metrics\_classification.py:1248: UndefinedMetricWarning: Precision a
nd F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` par
ameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
C:\pythonDSAI\lib\site-packages\sklearn\metrics\_classification.py:1248: UndefinedMetricWarning: Precision a
nd F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` par
ameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
C:\pythonDSAI\lib\site-packages\sklearn\metrics\_classification.py:1248: UndefinedMetricWarning: Precision a
nd F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` par
ameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
C:\pythonDSAI\lib\site-packages\sklearn\metrics\_classification.py:1248: UndefinedMetricWarning: Precision a
nd F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` par
ameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
C:\pythonDSAI\lib\site-packages\sklearn\metrics\_classification.py:1248: UndefinedMetricWarning: Precision a
nd F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` par
ameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
C:\pythonDSAI\lib\site-packages\sklearn\metrics\_classification.py:1248: UndefinedMetricWarning: Precision a
nd F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` par
ameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
C:\pythonDSAI\lib\site-packages\sklearn\metrics\_classification.py:1248: UndefinedMetricWarning: Precision a
nd F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` par
ameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
C:\pythonDSAI\lib\site-packages\sklearn\metrics\_classification.py:1248: UndefinedMetricWarning: Precision a
```

nd F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

```
_warn_prf(average, modifier, msg_start, len(result))
```

C:\pythonDSAI\lib\site-packages\sklearn\metrics_classification.py:1248: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

```
_warn_prf(average, modifier, msg_start, len(result))
```

C:\pythonDSAI\lib\site-packages\sklearn\metrics_classification.py:1248: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

	precision	recall	f1-score	support
-1.0	0.91	1.00	0.95	39
1.0	1.00	0.90	0.95	41
accuracy			0.95	80
macro avg	0.95	0.95	0.95	80
weighted avg	0.95	0.95	0.95	80

For kernel = gaussian
For C = 10

	precision	recall	f1-score	support
-1.0	0.91	1.00	0.95	39
1.0	1.00	0.90	0.95	41
accuracy			0.95	80
macro avg	0.95	0.95	0.95	80
weighted avg	0.95	0.95	0.95	80

In [46]:

```
from sklearn.svm import SVC

for kernel in ['linear', 'poly', 'rbf', 'sigmoid']:
    clf = SVC(C = 10, kernel = kernel)
    clf.fit(X_train, y_train.ravel())
    y_hat = clf.predict(X_test)
    print("For kernel = ", kernel)
    print(classification_report(y_test, y_hat))
```

For kernel = linear

	precision	recall	f1-score	support
-1.0	0.68	0.69	0.68	39
1.0	0.70	0.68	0.69	41
accuracy			0.69	80
macro avg	0.69	0.69	0.69	80
weighted avg	0.69	0.69	0.69	80

For kernel = poly

	precision	recall	f1-score	support
-1.0	0.89	0.62	0.73	39
1.0	0.72	0.93	0.81	41
accuracy			0.78	80
macro avg	0.80	0.77	0.77	80
weighted avg	0.80	0.78	0.77	80

For kernel = rbf

	precision	recall	f1-score	support
-1.0	0.76	0.82	0.79	39
1.0	0.82	0.76	0.78	41
accuracy			0.79	80
macro avg	0.79	0.79	0.79	80
weighted avg	0.79	0.79	0.79	80

For kernel = sigmoid

	precision	recall	f1-score	support
-1.0	0.41	0.49	0.45	39
1.0	0.41	0.34	0.37	41
accuracy			0.41	80
macro avg	0.41	0.41	0.41	80
weighted avg	0.41	0.41	0.41	80