

Problem Set 5 (Total points: 110 + bonus 50)

Support Vector Machines

In this problem set you will implement an SVM and fit it using quadratic programming. We will use the CVXOPT module to solve the optimization problems. You may want to start with solving the written problems at the end of this notebook or at least with reading the textbook. It will help a lot in some of the cells will take minutes to run, so feel free to test you code on smaller tasks while you go. Easiest way would be to remove both for-loops and run the code just once.

Quadratic Programming

The standard form of a QP can be formulated as

$$\begin{aligned} \min_x \quad & \frac{1}{2}x^T Px + q^T x \\ \text{subject to} \quad & Gx \leq h \\ & Ax = b \end{aligned}$$

where \leq is an element-wise \leq . CVXOPT solver finds an optimal solution x^* , given a set of matrices P, q, G, h, A, b .

FVL you can read on the methods to solve quadratic programming problems [here](#).

Problem 1. [10 points]

Design appropriate matrices to solve the following problem.

$$\begin{aligned} \min_x \quad & f(x) = x_1^2 + 4x_2^2 - 8x_1 - 16x_2 \\ \text{subject to} \quad & x_1 + x_2 \leq 5 \\ & x_1 \leq 3 \\ & x_2 \geq 0 \end{aligned}$$

```
In [102]: # run the code if the library is not installed
!pip install cvxopt
```

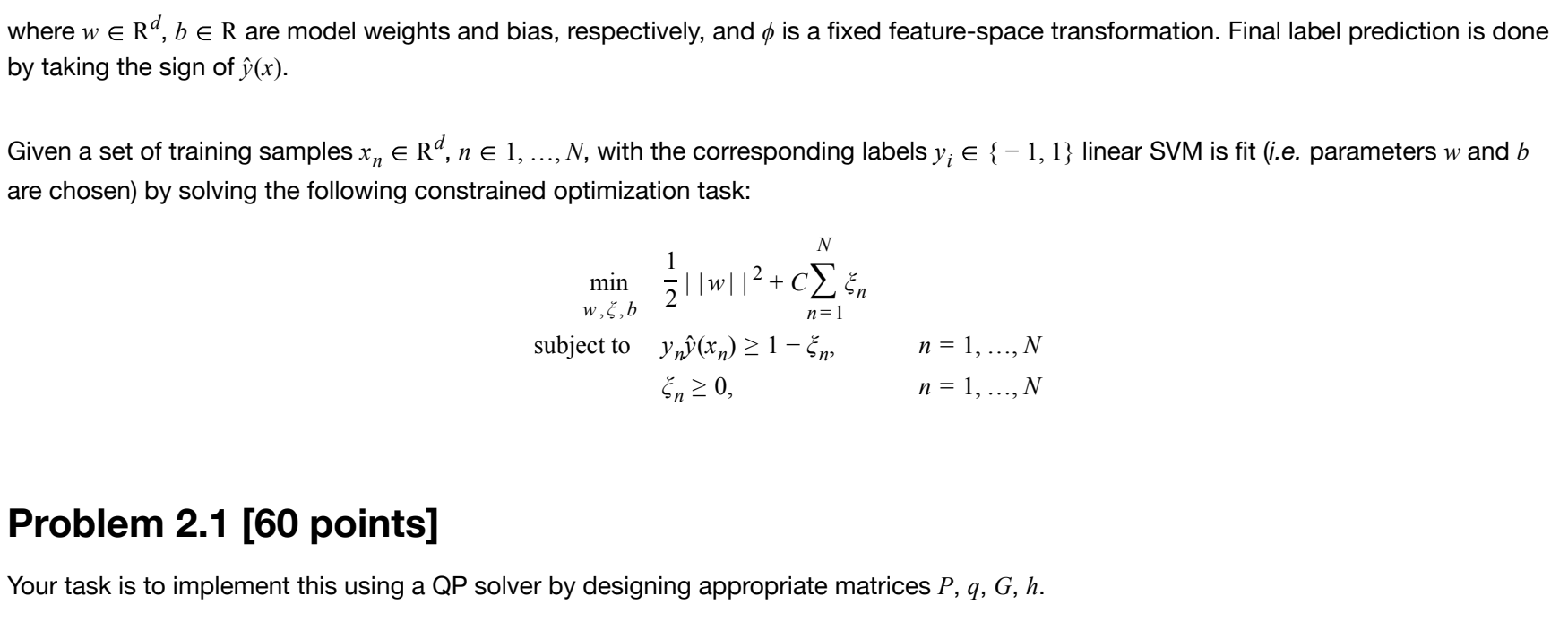
```
In [103]: from cvxopt import matrix, solvers
# Turns off the printing of CVXOPT solution for the rest of the notebook
solvers.options['show_progress'] = False

P = 2 * matrix([[1., 0.], [0., 4.]])
# Define q, G, h
q = matrix([-8., -16.])
G = matrix([[1., 1., 0.], [1., 0., -1.]])
h = matrix([[5., 3., 0.]])

sol = solvers.qp(P, q, G, h)
x1, x2 = sol['x']
print('Optimal x: ({:f}, {:f})'.format(x1, x2))

Optimal x: (2.99999993, 1.99927914)
```

Let's visualize the solution



Why is the solution not in the minimum? Because the x values have constraints

Linear SVM

Now, let's implement linear SVM. We will do this for a general case, that allows class distributions to overlap (see Bishop 7.1.1).

As a linear model, linear SVM produces classification scores for a given sample x as

$$\hat{y}(x) = w^T \phi(x) + b$$

where $w \in \mathbb{R}^d$, $b \in \mathbb{R}$ are model weights and bias, respectively, and ϕ is a fixed feature-space transformation. Final label prediction is done by taking the sign of $\hat{y}(x)$.

Given a set of training samples $x_n \in \mathbb{R}^d$, $y_n \in \{-1, 1\}$, with the corresponding labels $y_n \in \{-1, 1\}$ linear SVM is fit (i.e. parameters w and b are chosen) by solving the following constrained optimization task:

$$\begin{aligned} \min_{w, b, \xi} \quad & \frac{1}{2} \|w\|^2 + c \sum_{n=1}^N \xi_n \\ \text{subject to} \quad & y_n (w^T x_n + b) \geq 1 - \xi_n, \quad n = 1, \dots, N \\ & \xi_n \geq 0, \quad n = 1, \dots, N \end{aligned}$$

Problem 2.1 [60 points]

Your task is to implement this using a QP solver by designing appropriate matrices P, q, G, h .

Hints

1. You need to optimize over w, ξ, b . You can simply concatenate them into $z = (w, \xi, b)^T$ to feed it into QP-solver. Now, how to define the objective function and the constraints in terms of z ? (For example, $b_1 + b_2$ can be obtained from vector $(a_1, b_1, b_2, c_1, c_2)$ by taking the inner product with $(0, 1, 1, 0, 0)$.)
2. You can use `np.bmat` to construct matrices. Like this:

```
In [105]: np.bmat([[np.identity(3), np.zeros((3, 1))],
                [[np.zeros((2, 3)), -np.ones((2, 1))]])

Out[105]: matrix([[ 1.,  0.,  0.,  0.],
                  [ 0.,  1.,  0.,  0.],
                  [ 0.,  0.,  1.,  0.],
                  [ 0.,  0.,  0., -1.]])

In [113]: from sklearn.base import BaseEstimator
class LinearSVM(BaseEstimator):
    def __init__(self, C, transform=None):
        self.C = C
        self.transform = transform

    def fit(self, X, Y):
        """Fit Linear SVM using training dataset (X, Y).

        :param X: data samples of shape (N, d)
        :param Y: data target labels of size (N). Each label is either 1 or -1.

        # Apply transformation (phi) to X
        if self.transform is not None:
            X = self.transform(X)
            d = len(X[0])
            N = len(X)

        # Construct appropriate matrices here to solve the optimization problem described above.
        # We want optimal solution for vector (w, xi, b).

        # decision variables: w, xi, b -> N + d + 1 decision variables
        P = matrix(np.bmat([[np.identity(d), np.zeros((d, N + 1))], [np.zeros((N + 1, d + N + 1))]])
        q = matrix(np.bmat([[np.zeros(1, d)], self.C * np.ones(1, N)], np.zeros(1, 1)]))

        # C * sum(Xi)
        G = matrix(np.bmat([[np.zeros(1, d)], -np.multiply(X, Y.reshape((N, 1))), -np.identity(N), -Y.reshape((N, 1))],
                [[np.zeros((N, d), np.identity(N), np.zeros((N, 1))]])

        # h should be [-2, 0]
        h = matrix(np.bmat([[1 - np.ones(N, 1)], [np.zeros((N, 1))]])

        sol = solvers.qp(P, q, G, h)
        ans = np.array(sol['x']).flatten()
        self.weights_ = ans[1:d]
        self.bias_ = ans[d+1]

        # Find support vectors. Must be a boolean array of length N having True for support
        # vectors and False for the rest.
        self.support_vectors_ =

        # y_n = y * h(x_n) == 1 - xi_n
        margin = 1 - np.dot(X, self.weights_) + self.bias_
        andom_idx = np.argmax(margin, 1)
        self.support_vectors_ = np.take(self.support_vectors_, andom_idx)

    def predict_proba(self, X):
        """Make real-valued prediction for some new data.

        :param X: data samples of shape (N, d)
        :return: an array of N predicted scores.

        """
        if self.transform is not None:
            X = self.transform(X)
        return np.dot(self.weights_, X.T) + self.bias_

    def predict(self, X):
        """Make binary prediction for some new data.

        :param X: data samples of shape (N, d)
        :return: an array of N binary predicted labels from {-1, 1}.

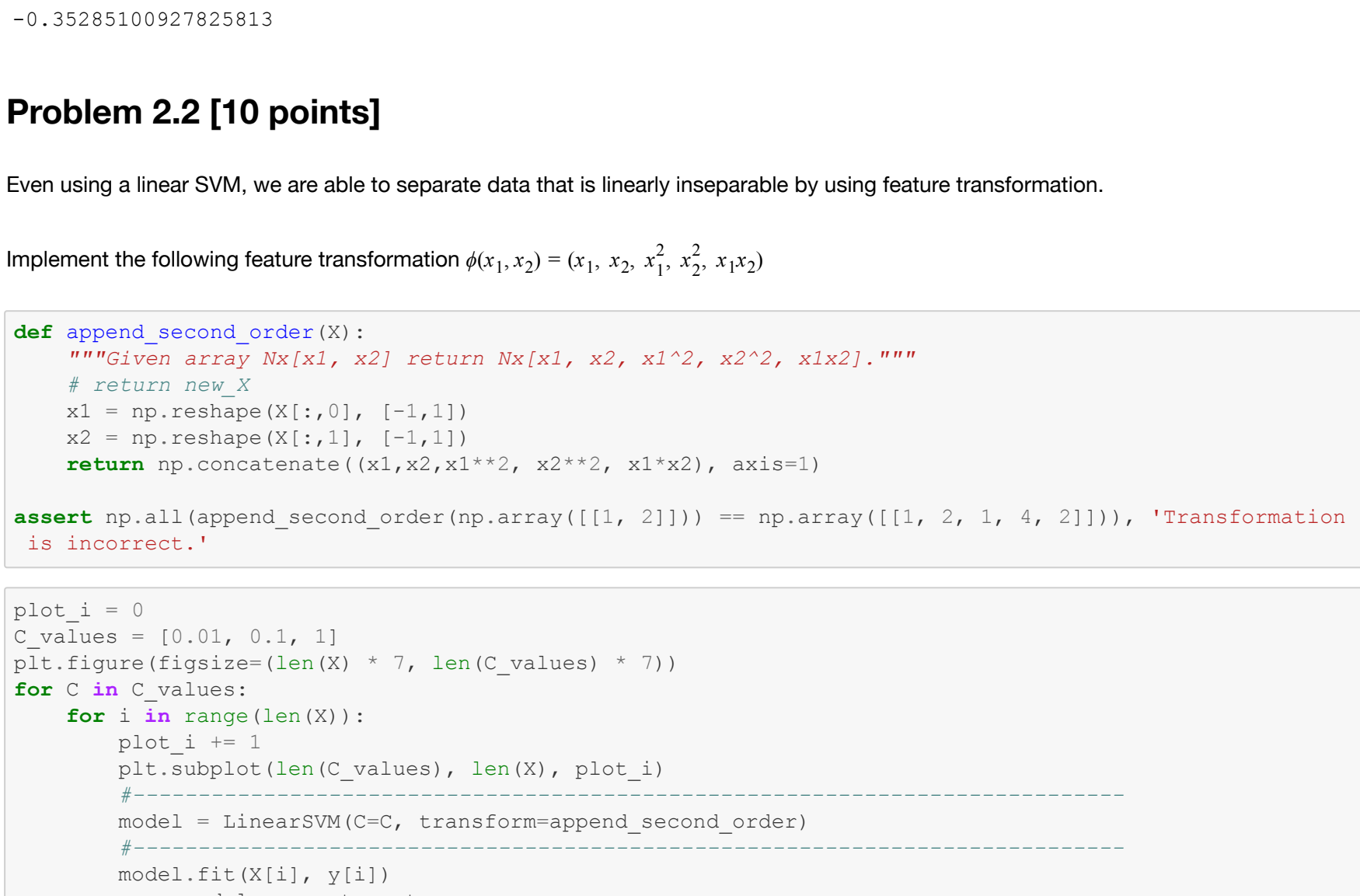
        """
        return np.sign(self.predict_proba(X))

Let's see how our LinearSVM performs on some data.
```

```
In [114]: from sklearn.datasets import make_classification, make_circles
X = [None, None, None]
Y = [None, None, None]
X[0], y[0] = make_classification(n_samples=100, n_features=2, n_redundant=0, n_clusters_per_class=1,
                                random_state=1)
X[1], y[1] = make_circles(n_samples=100, factor=0.5)
X[2], y[2] = make_classification(n_samples=100, n_features=2, n_redundant=0, n_clusters_per_class=1,
                                random_state=1)

# Go from [0, 1] to [-1, 1]
y = [2 * yy - 1 for yy in y]

In [115]: C_values = [0.01, 0.1, 1]
```



Why does the number of support vectors decrease as C increases? We are penalizing the slack variables more

For debug purposes. Very last model must have almost the same weights and bias:

$$w = \begin{pmatrix} -0.0784521 \\ 1.62264867 \end{pmatrix}$$
$$b = -0.35285100927825813$$

```
Out[116]: model.weights_
array([-0.0784521,  1.62264867])

In [117]: model.bias_

Out[117]: -0.35285100927825813
```

Problem 2.2 [10 points]

Even using a linear SVM, we are able to separate data that is linearly inseparable by using feature transformation.

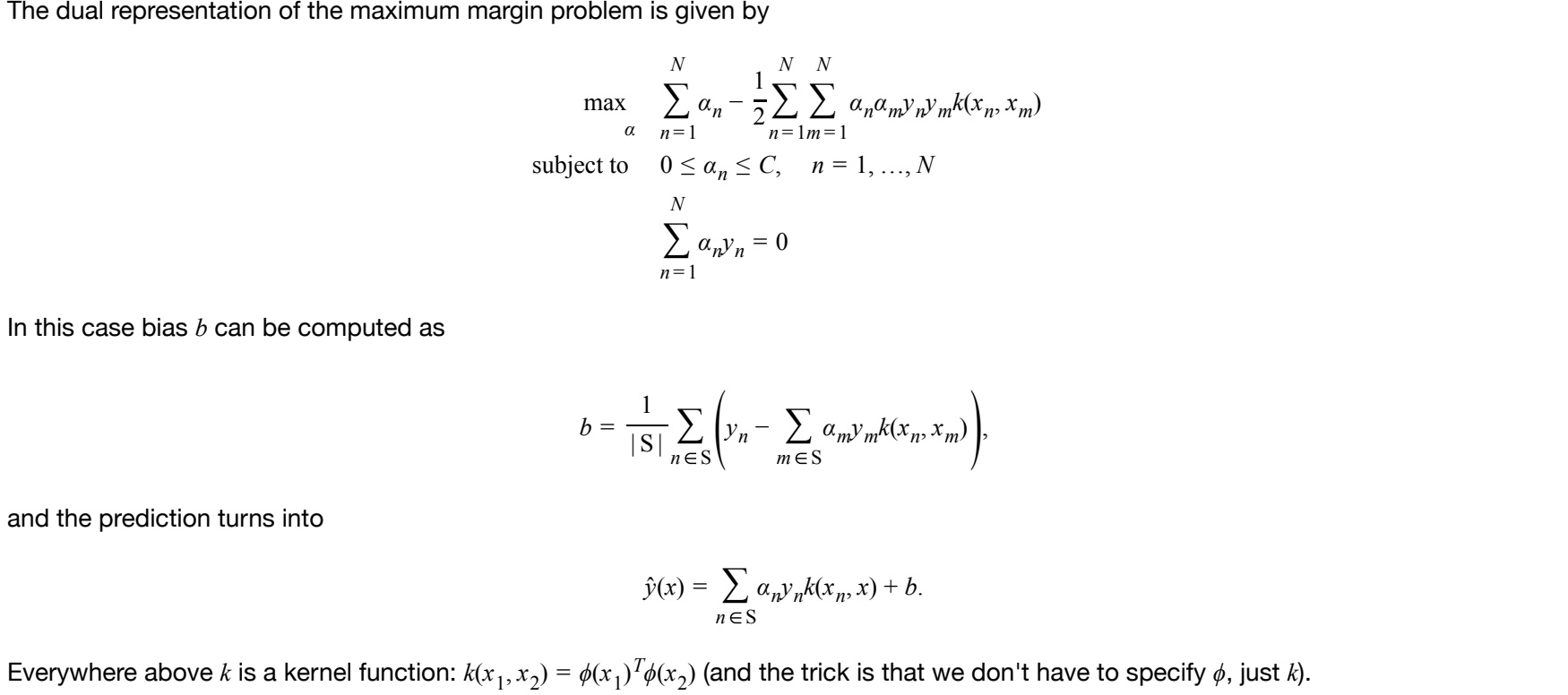
Implement the following feature transformation $\phi(x_1, x_2) = (x_1, x_2, x_1^2, x_2^2, x_1 x_2)^T$

```
In [126]: def append_second_order(X):
    """Given array Nx(x1, x2) return Nx(x1, x2, x1^2, x2^2, x1x2)."""
    # return new X
    x1 = np.reshape(X[:,0], [-1,1])
    x2 = np.reshape(X[:,1], [-1,1])
    xx, yy = np.meshgrid(x1, y2)
    return np.concatenate([X, x1**2, x2**2, x1*x2], axis=1)

assert np.all(append_second_order(np.array([[1, 2]])) == np.array([[1, 2, 1, 4, 2]])), "Transformation is incorrect"
```

```
In [127]: plot_i = 0
C_values = [0.01, 0.1, 1]
plt.figure(figsize=(len(X) * 7, len(C_values) * 7))
for C in C_values:
    for i in range(len(X)):
        plt.subplot(len(C_values), len(X), plot_i)
        #-----
        model = LinearSVM(C=C, transform=append_second_order)
        #-----
        model.fit(X[i], y[i])
        sv = model.support_vectors_
        n_sv = sv.sum()
        if n_sv > 0:
            plt.scatter(X[i][:, 0][sv], X[i][:, 1][sv], c=y[i][sv], cmap='autumn', marker='s',
                        linewidth=0.5, edgecolors=(0, 0, 1))
            if n_sv < len(X[i]):
                plt.scatter(X[i][:, 0][~sv], X[i][:, 1][~sv], c=y[i][~sv], cmap='autumn',
                            linewidth=0.5, edgecolors=(0, 0, 1))
            xvals = np.linspace(-3, 3, 200)
            yvals = np.linspace(-3, 3, 200)
            xx, yy = np.meshgrid(xvals, yvals)
            zz = np.reshape(model.predict_proba(np.c_[xx.ravel(), yy.ravel()]), (xx.shape))
            plt.pcolormesh(xx, yy, zz, cmap='autumn', zorder=0)
            plt.contour(xx, yy, zz, levels=(-1, 0, 1), colors='w', linewidths=1.5, zorder=1, linestyle='solid')

            plt.xlim([-3, 3])
            plt.ylim([-3, 3])
            plt.title('Dataset {i}, C={i}, #SV={i}'.format(i + 1, C, n_sv))
            plt.show()
```



Bonus part (Optional)

Dual representation. Kernel SVM

The dual representation of the maximum margin problem is given by

$$\begin{aligned} \max_{\alpha} \quad & \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m k(x_n, x_m) \\ \text{subject to} \quad & 0 \leq \alpha_n \leq C, \quad n = 1, \dots, N \\ & \sum_{n=1}^N \alpha_n b_n = 0 \end{aligned}$$

In this case bias b can be computed as

$$b = \frac{1}{|S|} \sum_{n \in S} \left(\sum_{m \in S} \alpha_n \alpha_m k(x_n, x_m) \right)$$

and the prediction turns into

$$\hat{y}(x) = \sum_{n=1}^N \alpha_n k(x_n, x) + b$$

Everywhere above k is a kernel function: $k(x_1, x_2) = \phi(x_1)^T \phi(x_2)$ (and the trick is that we don't have to specify ϕ , just k).

Note, that now

1. We want to maximize the objective function, not minimize it.
2. We have equality constraints. (That means we should use $=$ and b in qp-solver)
3. We need access to the support vectors (but not all the training samples) in order to make a prediction.

Problem 3.1 [40 points]

Implement KernelSVM

Hints

1. What is the variable α ?
2. How can we maximize a function given a tool for minimization?
3. What is the definition of a support vector in the dual representation?

```
In [ ]: class KernelSVM(BaseEstimator):
    def __init__(self, C, kernel=np.dot):
        self.C = C
        self.kernel = kernel

    def fit(self, X, Y):
        """Fit Kernel SVM using training dataset (X, Y).

        :param X: data samples of shape (N, d)
        :param Y: data target labels of size (N). Each label is either 1 or -1. Denoted as t_i in Bishop.

        """
        #-----
        # Construct appropriate matrices here to solve the optimization problem described above.
        # P =
        # q =
        # G =
        # h =
        #-----

        sol = solvers.qp(P, q, G, h, A, b)
        self.alpha_ = np.array(sol['x']).flatten()

        # Find support vectors. Must be a boolean array of length N having True for support
        # vectors and False for the rest.
        self.support_vectors_ =

        sv_ind = self.support_vectors_.nonzero()[0]
        self.X_sup = X[sv_ind]
        self.Y_sup = Y[sv_ind]
        self.alpha_sup = self.alpha[sv_ind]
        self.n_sup = len(sv_ind)

        #-----
        # Compute bias
        # self.bias_ =
        #-----

    def predict_proba(self, X):
        """Make real-valued prediction for some new data.

        :param X: data samples of shape (N, d)
        :return: an array of N predicted scores.

        """
        return y_hat

    def predict(self, X):
        """Make binary prediction for some new data.

        :param X: data samples of shape (N, d)
        :return: an array of N binary predicted labels from {-1, 1}.

        """
        return np.sign(self.predict_proba(X))

We can first test our implementation by using the dot product as a kernel function. What should we expect in this case?
```

```
In [ ]: C = 1
d = 10
plt.figure(figsize=(7, 7))
#-----
model = KernelSVM(C=C, kernel=np.dot)
#-----
model.fit(X[1], y[1])
sv = model.support_vectors_
n_sv = sv.sum()
if n_sv > 0:
    plt.scatter(X[1][:, 0][sv], X[1][:, 1][sv], c=y[1][sv], cmap='autumn', marker='s',
                linewidth=0.5, edgecolors=(0, 0, 1))
    if n_sv < len(X[1]):
        plt.scatter(X[1][:, 0][~sv], X[1][:, 1][~sv], c=y[1][~sv], cmap='autumn',
                    linewidth=0.5, edgecolors=(0, 0, 1))
    xvals = np.linspace(-3, 3, 200)
    yvals = np.linspace(-3, 3, 200)
    xx, yy = np.meshgrid(xvals, yvals)
    zz = np.reshape(model.predict_proba(np.c_[xx.ravel(), yy.ravel()]), (xx.shape))
    plt.pcolormesh(xx, yy, zz, cmap='autumn', zorder=0)
    plt.contour(xx, yy, zz, levels=(-1, 0, 1), colors='w', linewidth=1.5, zorder=1, linestyle='solid')

    plt.xlim([-3, 3])
    plt.ylim([-3, 3])
    plt.title('Dataset {i}, C={i}, d={i}, #SV={i}'.format(i + 1, C, d, n_sv))
    plt.show()
```

Problem 3.2 [5 points]

Implement a polynomial kernel function ([link](#)).

```
In [ ]: def polynomial_kernel(d, c=0):
    """Returns a polynomial kernel FUNCTION."""
    def kernel(x, y):
        """
        :param x: vector of size L
        :param y: vector of size L
        :return: (polynomial kernel of degree d with bias parameter c) of x and y. A scalar.

        """
        pass
    return kernel

assert polynomial_kernel(d=2, c=1)(np.array([1, 2]), np.array([3, 4])) == 144, "Polynomial kernel implemented incorrectly"
```

Let's see how it performs. This might take some time to run.

```
In [ ]: plot_i = 0
C = 10
d_values = [2, 3, 4]
plt.figure(figsize=(len(X) * 7, len(d_values) * 7))
for d in d_values:
    for i in range(len(X)):
        plt.subplot(len(d_values), len(X), plot_i)
        #-----
        model = KernelSVM(C=C, kernel=polynomial_kernel(d))
        #-----
        model.fit(X[i], y[i])
        sv = model.support_vectors_
        n_sv = sv.sum()
        if n_sv > 0:
            plt.scatter(X[i][:, 0][sv], X[i][:, 1][sv], c=y[i][sv], cmap='autumn', marker='s',
                        linewidth=0.5, edgecolors=(0, 0, 1))
            if n_sv < len(X[i]):
                plt.scatter(X[i][:, 0][~sv], X[i][:, 1][~sv], c=y[i][~sv], cmap='autumn',
                            linewidth=0.5, edgecolors=(0, 0, 1))
            xvals = np.linspace(-3, 3, 200)
            yvals = np.linspace(-3, 3, 200)
            xx, yy = np.meshgrid(xvals, yvals)
            zz = np.reshape(model.predict_proba(np.c_[xx.ravel(), yy.ravel()]), (xx.shape))
            plt.pcolormesh(xx, yy, zz, cmap='autumn', zorder=0)
            plt.contour(xx, yy, zz, levels=(-1, 0, 1), colors='w', linewidth=1.5, zorder=1, linestyle='solid')

            plt.xlim([-3, 3])
            plt.ylim([-3, 3])
            plt.title('Dataset {i}, C={i}, d={i}, #SV={i}'.format(i + 1, C, d, len(model.support_vectors_)))
            plt.show()
```

Task 3.3 [5 points]

Finally, you need to implement a radial basis function kernel ([link](#)).

```
In [ ]: def RBF_kernel(sigma):
    """Returns an RBF kernel FUNCTION."""
    def kernel(x, y):
        """
        :param x: vector of size L
        :param y: vector of size L
        :return: (rbf kernel with parameter sigma) of x and y. A scalar.

        """
        pass
    return kernel

Let's see how it performs. This might take some time to run.
```

```
In [ ]: plot_i = 0
C = 10
sigma_values = [0.1, 1, 10]
plt.figure(figsize=(len(sigma_values) * 7, len(C_values) * 7))
for C in C_values:
    for sigma in sigma_values:
        plot_i += 1
        plt.subplot(len(C_values), len(X), plot_i)
        #-----
        model = KernelSVM(C=C, kernel=RBF_kernel(sigma))
        #-----
        model.fit(X[1], y[1])
        sv = model.support_vectors_
        n_sv = sv.sum()
        if n_sv > 0:
            plt.scatter(X[1][:, 0][sv], X[1][:, 1][sv], c=y[1][sv], cmap='autumn', marker='s',
                        linewidth=0.5, edgecolors=(0, 0, 1))
            if n_sv < len(X[1]):
                plt.scatter(X[1][:, 0][~sv], X[1][:, 1][~sv], c=y[1][~sv], cmap='autumn',
                            linewidth=0.5, edgecolors=(0, 0, 1))
            xvals = np.linspace(-3, 3, 200)
            yvals = np.linspace(-3, 3, 200)
            xx, yy = np.meshgrid(xvals, yvals)
            zz = np.reshape(model.predict_proba(np.c_[xx.ravel(), yy.ravel()]), (xx.shape))
            plt.pcolormesh(xx, yy, zz, cmap='autumn', zorder=0)
            plt.contour(xx, yy, zz, levels=(-1, 0, 1), colors='w', linewidth=1.5, zorder=1, linestyle='solid')

            plt.xlim([-3, 3])
            plt.ylim([-3, 3])
            plt.title('Dataset {i}, C={i}, sigma={i}, #SV={i}'.format(i + 1, C, sigma, n_sv))
            plt.show()
```

Well done!

Awesome! Now you understand all of the important parameters in SVMs. Have a look at SVM from scikit-learn module and how it is used (very similar to ours).

```
In [ ]: from sklearn.svm import SVC
SVC()
```

```
In [ ]: plot_i = 0
C = 10
d_values = [2, 3, 4]
plt.figure(figsize=(len(X) * 7, len(d_values) * 7))
for d in d_values:
    for i in range(len(X)):
        plt.subplot(len(d_values), len(X), plot_i)
        #-----
        model = SVC(kernel='poly', degree=d, gamma='auto', probability=True)
        #-----
        model.fit(X[1], y[1])
        sv = model.support_vectors_
        n_sv = sv.sum()
        if n_sv > 0:
            plt.scatter(X[1][:, 0][sv], X[1][:, 1][sv], c=y[1][sv], cmap='autumn', marker='s',
                        linewidth=0.5, edgecolors=(0, 0, 1))
            if n_sv < len(X[1]):
                plt.scatter(X[1][:, 0][~sv], X[1][:, 1][~sv], c=y[1][~sv], cmap='autumn',
                            linewidth=0.5, edgecolors=(0, 0, 1))
            xvals = np.linspace(-3, 3, 200)
            yvals = np.linspace(-3, 3, 200)
            xx, yy = np.meshgrid(xvals, yvals)
            zz = np.reshape(model.predict_proba(np.c_[xx.ravel(), yy.ravel()]), (xx.shape))
            plt.pcolormesh(xx, yy, zz, cmap='autumn', zorder=0)
            plt.contour(xx, yy, zz, levels=(-1, 0, 1), colors='w', linewidth=1.5, zorder=1, linestyle='solid')

            plt.xlim([-3, 3])
            plt.ylim([-3, 3])
            plt.title('Dataset {i}, C={i}, d={i}, #SV={i}'.format(i + 1, C, d, len(model.support_vectors_)))
            plt.show()
```

4 Written Problems

Problem 4.1 Dual Representations [10 pts]

Read Section 6.1 in Bishop, and work through all of the steps of the derivations in equations 6.2-6.9. You should understand how the derivation works in detail. Write down your understanding.

$$J(W) = \frac{1}{2} \sum_{n=1}^N (W^T \phi(X_n) - t_n)^2 + \frac{\lambda}{2} W^T W \quad (6.2) \quad \frac{\partial J}{\partial W} = \frac{\partial}{\partial W} \left(\frac{1}{2} \sum_{n=1}^N (W^T \phi(X_n) - t_n)^2 + \lambda W^T W \right)$$

solve for W :

$$\sum_{n=1}^N W^T \phi(X_n) - t_n + \phi(X_n)^T W = 0$$
$$\sum_{n=1}^N W^T \phi(X_n) - t_n + \phi(X_n)^T W = -\lambda W$$
$$W = \frac{-1}{\lambda} \sum_{n=1}^N W^T \phi(X_n) - t_n$$

define

$$a_n = \frac{1}{\lambda} (W^T \phi(X_n) - t_n) \quad (6.4)$$

then

$$W = \sum_{n=1}^N a_n \phi(X_n) = \Phi^T a \quad \text{where } \Phi \text{ is the design matrix.} \quad (6.3)$$

substitute $W = \Phi^T a$ into $J(W)$:

$$J(a) = \frac{1}{2} (a^T \Phi - t)^T (a^T \Phi - t) + \frac{\lambda}{2} a^T \Phi \Phi^T a$$
$$J(a) = \frac{1}{2} (a^T \Phi \Phi^T a - 2a^T \Phi^T t + t^T t) + \frac{\lambda}{2} a^T \Phi \Phi^T a \quad (6.5)$$

define $K = \Phi \Phi^T$ where

$$K_{nm} = \phi(X_n)^T \phi(X_m) = K(X_n, X_m) \quad (6.6)$$

substituting K into J

$$J(a) = \frac{1}{2} a^T K a - a^T K t + \frac{1}{2} t^T t + \frac{\lambda}{2} a^T K a \quad (6.7)$$

take derivative with respect to a and then solve for a :

$$\frac{\partial J}{\partial a} = K K a - K t + \lambda K a = 0$$
$$a K K + \lambda K a = K t$$
$$a (K K + \lambda K) = K t$$
$$a (K + \lambda I_N)^{-1} = K t \quad (6.8)$$

substitute back into the linear regression model:

$$y(X) = W^T \phi(X) = a^T \Phi \phi(X) = K(X)^T (K + \lambda I_N)^{-1} t \quad (6.9)$$

Problem 4.2 Kernels [10 pts]

Read Section 6.2 and Verify the results (6.13) and (6.14) for constructing valid kernels.

(6.13)

C is a constant. Since it is positive, it meets the constraint. Then $k(x, x') = c \delta(x, x')$ holds because $k(x, x')$ is a valid kernel.

We can just change the transformation function ϕ to account for c .

(6.14)

Since ℓ_1 is any function, we change the transformation function linearly to account for ℓ_1 . Then, since $k_1(x, x') = \ell_1(x, x') + \ell_2(x, x')$ and $k_2(x, x') = \ell_2(x, x')$ is also valid.

Problem 4.3 Maximum Margin Classifiers [10 pts]

Read section 7.1 and show that, if the 1 on the right hand side of the constraint (7.5) is replaced by some arbitrary constant $\gamma > 0$, the solution for maximum margin hyperplane is unchanged.

$$t_d(\mathbf{w}^T \phi(x_d) + b) \geq 1$$

Let γ be any arbitrary constant greater than 0. Then

$$t_d(W^T \phi(x_d) + b) \geq \gamma$$

$$t_d(\frac{W^T \phi(x_d)}{\gamma} + \frac{b}{\gamma}) \geq 1$$

Scaling by a constant does not change the hyperplane so the solution is the same.

In []: