# Support Vector Machines

In this exercise sheet, you will experiment with training various support vector machines on a subset of the MNIST dataset composed of digits 5 and 6. First, download the MNIST dataset from http://yann.lecun.com/exdb/mnist/, uncompress the downloaded files, and place them in a `data/` subfolder. Install the optimization library CVXOPT (`python-cvxopt` package, or directly from the website `www.cvxopt.org`). This library will be used to optimize the dual SVM in part A.

## Part A: Kernel SVM and Optimization in the Dual

We would like to learn a nonlinear SVM by optimizing its dual. An advantage of the dual SVM compared to the primal SVM is that it allows to use nonlinear kernels such as the Gaussian kernel, that we define as:

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

The dual SVM consists of solving the following quadratic program:

$$\max_{\alpha} \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j k(x_i, x_j)$$

subject to:

$$0 \leq \alpha_i \leq C \qquad \text{and} \qquad \sum_{i=1}^{N} \alpha_i y_i = 0.$$

Then, given the alphas, the prediction of the SVM can be obtained as:

$$f(x) = \begin{cases} 1 & \text{if} \quad \sum_{i=1}^{N} \alpha_i y_i k(x, x_i) + \theta > 0 \\ -1 & \text{if} \quad \sum_{i=1}^{N} \alpha_i y_i k(x, x_i) + \theta < 0 \end{cases}$$

where

$$\theta = \frac{1}{\#SV} \sum_{i \in SV} \left( y_i - \sum_{j=1}^{N} \alpha_j y_j k(x_i, x_j) \right)$$

and `SV` is the set of indices corresponding to the unbound support vectors.

### Implementation (25 P)

We will solve the dual SVM applied to the MNIST dataset using the CVXOPT quadratic optimizer. For this, we have to build the data structures (vectors and matrices) to must be passed to the optimizer.

- *Implement* a function `gaussianKernel` that returns for a Gaussian kernel of scale $\sigma$, the Gram matrix of the two data sets given as argument.

- *Implement* a function `getQPMatrices` that builds the matrices P, q, G, h, A, b (of type `cvxopt.matrix`) that need to be passed as argument to the optimizer `cvxopt.solvers.qp`.

- *Run* the code below using the functions that you just implemented. (It should take less than 3 minutes.)

```
In [1]: import utils,numpy,cvxopt,cvxopt.solvers
        import solutions

        Xtrain,Ttrain,Xtest,Ttest = utils.getMNIST56()
```

```
cvxopt.solvers.options['show_progress'] = False

for scale in [10,30,100]:
    for C in [1,10,100]:

        # Prepare kernel matrices

        ### TODO: REPLACE BY YOUR OWN CODE
        Ktrain = solutions.gaussianKernel(Xtrain,Xtrain,scale)
        Ktest  = solutions.gaussianKernel(Xtest,Xtrain,scale)
        ###


        # Prepare the matrices for the quadratic program

        ### TODO: REPLACE BY YOUR OWN CODE
        P,q,G,h,A,b = solutions.getQPMatrices(Ktrain,Ttrain,C)
        ###


        # Train the model (i.e. compute the alphas)
        alpha = numpy.array(cvxopt.solvers.qp(P,q,G,h,A,b)['x']).flatten()

        # Get predictions for the training and test set
        SV = (alpha>1e-6)
        uSV = SV*(alpha<C-1e-6)
        theta = 1.0/sum(uSV)*(Ttrain[uSV]-numpy.dot(Ktrain[uSV,:],alpha*Ttrain)).sum()
        Ytrain = numpy.sign(numpy.dot(Ktrain[:,SV],alpha[SV]*Ttrain[SV])+theta)
        Ytest  = numpy.sign(numpy.dot(Ktest [:,SV],alpha[SV]*Ttrain[SV])+theta)

        # Print accuracy and number of support vectors
        Atrain = (Ytrain==Ttrain).mean()
        Atest  = (Ytest ==Ttest ).mean()
        print('Scale=%3d  C=%3d  SV: %4d  Train: %.3f  Test: %.3f'%(scale,C,sum(SV),Atrain,Atest
    print('')

Scale= 10  C=  1  SV: 1000  Train: 1.000  Test: 0.937
Scale= 10  C= 10  SV: 1000  Train: 1.000  Test: 0.937
Scale= 10  C=100  SV: 1000  Train: 1.000  Test: 0.937

Scale= 30  C=  1  SV:  254  Train: 1.000  Test: 0.985
Scale= 30  C= 10  SV:  274  Train: 1.000  Test: 0.986
Scale= 30  C=100  SV:  256  Train: 1.000  Test: 0.986

Scale=100  C=  1  SV:  317  Train: 0.973  Test: 0.971
Scale=100  C= 10  SV:  159  Train: 0.990  Test: 0.975
Scale=100  C=100  SV:  136  Train: 1.000  Test: 0.975
```

**Analysis (10 P)**

- *Explain* which combinations of parameters $\sigma$ and $C$ lead to good generalization, underfitting or overfitting?

- *Explain* which combinations of parameters $\sigma$ and $C$ produce the fastest classifiers (in terms of amount of computation needed at prediction time)?

## Part B: Linear SVMs and Gradient Descent in the Primal

The quadratic problem of the dual SVM does not scale well with the number of data points. For large number of data points, it is generally more appropriate to optimize the SVM in the primal. The primal optimization problem for linear SVMs can be written as

$$\min_{w,\theta} ||w||^2 + C \sum_{i=1}^{N} \xi_i \qquad \text{where} \qquad \forall_{i=1}^{N} : y_i(w \cdot x_i + \theta) \geq 1 - \xi_i \qquad \text{and} \qquad \xi_i \geq 0.$$

It is common to incorporate the constraints directly into the objective and then minimizing the unconstrained objective

$$J(w,\theta) = ||w||^2 + C \sum_{i=1}^{N} \max(0, 1 - y_i(w \cdot x_i + \theta))$$

using simple gradient descent.

### Implementation (15 P)

- *Implement* the function J computing the objective $J(w,\theta)$

- *Implement* the function DJ computing the gradient of the objective $J(w,\theta)$ with respect to the parameters $w$ and $\theta$.

- *Run* the code below using the functions that you just implemented. (It should take less than 1 minute.)

```
In [2]: import utils,numpy
        import solutions

        C = 10.0
        lr = 0.001

        Xtrain,Ttrain,Xtest,Ttest = utils.getMNIST56()

        n,d = Xtrain.shape

        w = numpy.zeros([d])
        theta = 1e-9

        for it in range(0,101):

            # Monitor the training and test error every 5 iterations
            if it%5==0:
                Ytrain = numpy.sign(numpy.dot(Xtrain,w)+theta)
                Ytest  = numpy.sign(numpy.dot(Xtest ,w)+theta)

                ### TODO: REPLACE BY YOUR OWN CODE
                Obj    = solutions.J(w,theta,C,Xtrain,Ttrain)
                ###

                Etrain = (Ytrain==Ttrain).mean()
                Etest  = (Ytest ==Ttest ).mean()
                print('It=%3d   J: %9.3f  Train: %.3f  Test: %.3f'%(it,Obj,Etrain,Etest))

            ### TODO: REPLACE BY YOUR OWN CODE
            dw,dtheta = solutions.DJ(w,theta,C,Xtrain,Ttrain)
            ###
```

```
            w = w - lr*dw
            theta = theta - lr*dtheta

It=  0   J: 10000.000  Train: 0.471  Test: 0.482
It=  5   J: 68520.417  Train: 0.961  Test: 0.958
It= 10   J: 49918.674  Train: 0.973  Test: 0.961
It= 15   J: 37473.229  Train: 0.973  Test: 0.963
It= 20   J: 28590.129  Train: 0.974  Test: 0.965
It= 25   J: 21746.877  Train: 0.977  Test: 0.967
It= 30   J: 16987.200  Train: 0.980  Test: 0.968
It= 35   J: 13646.095  Train: 0.986  Test: 0.967
It= 40   J: 11187.127  Train: 0.986  Test: 0.967
It= 45   J:  9182.940  Train: 0.991  Test: 0.967
It= 50   J:  7692.273  Train: 0.990  Test: 0.968
It= 55   J:  6437.609  Train: 0.988  Test: 0.966
It= 60   J:  5253.071  Train: 0.995  Test: 0.966
It= 65   J:  4515.520  Train: 0.992  Test: 0.967
It= 70   J:  4016.851  Train: 0.996  Test: 0.966
It= 75   J:  3647.983  Train: 0.997  Test: 0.965
It= 80   J:  3497.204  Train: 0.998  Test: 0.966
It= 85   J:  3404.280  Train: 1.000  Test: 0.966
It= 90   J:  3336.804  Train: 1.000  Test: 0.966
It= 95   J:  3270.665  Train: 1.000  Test: 0.966
It=100   J:  3205.837  Train: 1.000  Test: 0.966
```