

# DD2431 Machine Learning

## Lab 4: Support Vector Machines

November 26, 2008

### 1 Background

Support vector machines, when used for classification, find a hyperplane  $\mathbf{w}^T \mathbf{x} + b = 0$  that separates two classes of data. Here  $\mathbf{w}$  and  $b$  are parameters specifying the hyperplane and  $\mathbf{x}$  are points in space. If the data to be classified is written as a set of points  $\mathbf{x}_i$  along with class labels  $y_i \in \{-1, +1\}$  the condition that the data is correctly classified becomes  $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 0$ .

#### 1.1 Primal formulation

Of many potential hyperplanes that might fulfil this basic separation criterion, we would like to select one such that *the distance between the hyperplane and any data point is as large as possible*. This distance is known as the *margin*. When the hyperplane is put on a canonical form (more than one parameter set  $\mathbf{w}, b$  may describe the same hyperplane) this gives the following optimisation problem. Since we will later derive a “dual” version, we refer to this as the *primal* problem.

$$\text{minimise:} \quad \mathbf{w}^T \mathbf{w} \quad (1)$$

$$\text{when:} \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad i \in \{1..m\} \quad (2)$$

The data points  $\mathbf{x}_i$  closest to the separating hyperplane will be precisely the ones where equation 2 is satisfied with equality. Those data points are known as *support vectors*. In the non-degenerate case they will be the only data points directly determining the location of the separating hyperplane, in the sense that a small displacement of these points would shift the separating hyperplane. The resulting hard classifier is:

$$f(x) = \text{sgn}(\mathbf{w}^T \mathbf{x} + b) \quad (3)$$

## 1.2 Dual formulation

Equivalently, we may write the *dual* problem. We arrive at the dual problem by assigning positive Lagrange multipliers  $\alpha_i$  to the inequality constraints in equation 2. The dual problem becomes:

$$\text{maximise: } \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \quad (4)$$

$$\text{when: } \sum_{i=1}^m \alpha_i y_i = 0 \quad (5)$$

$$0 \leq \alpha_i \quad i \in \{1..m\} \quad (6)$$

We may recognise the support vectors also in the dual formulation. They are the data points for which  $\alpha_i > 0$ . Once the dual problem has been solved, the primary variables are retrieved as

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \quad (7)$$

$$b = -\frac{1}{2} \mathbf{w}^T (\mathbf{x}_p + \mathbf{x}_n) \quad (8)$$

where  $\mathbf{x}_p$  is a positive support vector, having  $\alpha_p > 0$  and  $y_p = +1$  and  $\mathbf{x}_n$  is a negative support vector, having  $\alpha_n > 0$  and  $y_n = -1$ . Alternatively, for better numerical precision, all support vectors may be used to calculate the bias;

$$b = \frac{1}{|\{i : \alpha_i \neq 0\}|} \sum_{i: \alpha_i \neq 0} (y_i - \mathbf{w}^T \mathbf{x}_i) \quad (9)$$

where the denominator in the first factor is the count of support vectors.

## 2 Build your own SVM

In this part of the lab, you will build your own support vector machine, by directly solving the above optimisation problems. You will be doing this twice, once using the primal formulation, once using the dual formulation.

The data you will be working on is linearly separable, available in `Matlab` format in the file `linsep.mat`. The data file contains two classes. The positive instances are listed in the data structure `classA`, the negative instances in `classB`. For your convenience, there are also ready-made data structures `x` and `y`. They contain all of the data points and class labels, respectively. The vector `permute` describes how they were created from `classA` and `classB`.

## 2.1 Primal formulation

You will perform a constrained optimisation in the variables  $[\mathbf{w}, b]$ . You may use the Matlab function `fmincon` for this. This function takes as arguments an objective function and the applicable constraints. You should not have any problem passing the objective function  $\mathbf{w}^T \mathbf{w}$  to `fmincon`. Just define it in a separate `myfun.m` file and pass it as `@myfun`. Remember that the elements of the vector  $\mathbf{w}$  make up the first part of the variable  $\mathbf{X}$ , in which you are optimising, the bias  $b$  is the final element. The constraints (equation 2) are linear. Therefore, they may be expressed as a matrix  $\mathbf{A}$  and a vector  $\mathbf{B}$ ; the second and third parameters to `fmincon`. (Hint: The vector  $\mathbf{B}$  can be made to have all its elements equal  $-1$ .) When the optimisation has finished, you may find which inequality constraints are active in the return value `LAMBDA.ineqlin`.

**Task:** Plot the data points for the two classes along with the optimal separating hyperplane. Mark the support vectors in the plot, e.g. by circles; `plot(x,y,'o')`. Make sure that your plot is readable in black and white. You may use the below code to do some of the plotting.

**Question:** What is the dimension of the  $\mathbf{w}$  vector?

**Question:** How many support vectors do you get?

In the code below, `myfunction` is your classifier function, which is equation 3, *without* the thresholding “`sgn`” part. It should accept an  $\mathbf{x}$  vector as its first argument. Other arguments that your function requires, such as the  $\mathbf{w}$  vector and the bias  $b$ , can be listed last in the call to `meshapply`, as indicated. The function file `meshapply.m` is found in the lab directory.

```
clf; hold on;

[X,Y] = meshgrid(min(x(:,1)):.1:max(x(:,1)), ...
                 min(x(:,2)):.1:max(x(:,2)));

Z = meshapply(@myfunction, X, Y, [other args]);

contour(X, Y, Z, [-1 0 1]);
plot (classA(:,1), classA(:,2), 'r*')
plot (classB(:,1), classB(:,2), 'b*')
```

## 2.2 Dual formulation

This time, the optimisation happens over the Lagrangian vector  $\alpha$ . What is its dimension? Again, the `Matlab` function `fmincon` may be useful to solve the problem. If you would like, you can calculate the matrix  $K_{ij} := \mathbf{x}_i^T \mathbf{x}_j$  (or even  $M_{ij} := y_i y_j K_{ij}$ ) in advance and have the objective function `FUN` access  $K$  or  $M$  as a global variable. Note that you need linear equality constraints in addition to inequality constraints. You can accomplish this by using the parameters `Aeq` and `Beq` of `fmincon`.

**Task:** Again plot the data points, the separating hyperplane and mark the support vectors.

**Question:** How many support vectors do you get this time?

**Question:** Is the separating hyperplane different from what you got using the primal formulation? Can you see from looking at the geometry of the problem (the data points near the border between classes) if there is a unique maximal margin hyperplane in this case?

## 3 Linearly non-separable data

The above methods only work for linearly separable data.

**Task:** Redo the above experiments using the data set `nlinsep.mat`. Do not attempt to fix your code, just run it naively for the new data.

**Question:** What happens when you try to apply the primal and dual approaches, respectively?

When data are not linearly separable, some data points will necessarily be misclassified. We loosen the bounds in equation 2 to allow for this, but any deviation is penalised in a new objective function (equation 10). The primal problem now reads

$$\text{minimise: } \quad \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^m \xi_i^\mu \quad (10)$$

$$\text{when: } \quad y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad i \in \{1..m\} \quad (11)$$

$$\xi_i \geq 0 \quad (12)$$

The  $\xi_i$  measure the degree of misclassification for each data point. Note that they are new variables in the problem. They will be zero for data points that

are correctly classified and outside of the margin.  $C$  (a coefficient) and  $\mu$  (an exponent) are positive parameters that determine how misclassifications are penalised.

The dual formulation becomes particularly simple when the parameter  $\mu = 1$ , so for simplicity we will constrain ourselves to this case for the lab. When  $\mu = 1$ , the dual problem is the same as before, except for the constraints on the dual variables. These are now also bounded from above:

$$0 \leq \alpha_i \leq C \quad i \in \{1..m\} \quad (13)$$

In addition, we have to restrict the set of data points over which we compute the bias to those with  $\xi_i = 0$ , meaning  $\alpha_i \leq C$ <sup>1</sup>. We replace equation 9 with the following:

$$b = \frac{1}{|\{i : 0 < \alpha_i < C\}|} \sum_{i: 0 < \alpha_i < C} (y_i - \mathbf{w}^T \mathbf{x}_i) \quad (14)$$

### 3.1 Classify linearly non-separable data

Modify the *dual* version of your algorithm according to equation 13. You will have to introduce the new constraints involving  $C$  from equation 13 and modify the way the bias is calculated according to equation 14.

**Task:** Classify the linearly non-separable data using your modified algorithm. Plot the data points and the separating hyperplane as before. Based on your visual impression of how the separating hyperplane is positioned, find your favourite value for the parameter  $C$ .

**Question (optional):** How does changing the value of  $C$  alter the placement of the separating hyperplane and the number of support vectors? Show a representative result. (It may be hard to see what is happening to the hyperplane; it will be easier in the next section.)

## 4 The kernel trick

It is often advantageous to map data into a higher-dimensional space, since a problem that was not linearly separable in the original space, can often

---

<sup>1</sup>If there are no data points  $0 \leq \alpha_i \leq C$  we could in principle still compute the bias from equation 11, but for this lab we will just give up in that case.

be made so in higher dimensional space. The real power of support vector machines comes from the following realisation: If the data points are mapped to a different space;  $x \mapsto \phi(\mathbf{x})$  before finding a hyperplane to separate them, the only operation that needs to be carried out in the new space is the computation of the inner product  $\phi^T(\mathbf{x}_1)\phi(\mathbf{x}_2)$ .

We denote the inner product in the new space  $\phi(\mathfrak{R}^m)$  as  $K_\phi(\mathbf{x}_1, \mathbf{x}_2) := \phi^T(\mathbf{x}_1)\phi(\mathbf{x}_2)$  and refer to it as a *kernel function*. For every transformation  $\phi(\mathbf{x})$  there exists, by this definition, a kernel  $K_\phi(\mathbf{x}_1, \mathbf{x}_2)$ .

**Question (optional):** Not every function  $f(\mathbf{x}_1, \mathbf{x}_2) : \mathfrak{R}^m \times \mathfrak{R}^m \mapsto \mathfrak{R}$  is a kernel. There are many functions that cannot be expressed as  $f(\mathbf{x}_1, \mathbf{x}_2) = \phi^T(\mathbf{x}_1)\phi(\mathbf{x}_2)$  for *any* transformation  $\phi(\mathbf{x})$ . Can you find an example of a function that cannot be?

## 4.1 Disposing of $w$

We now proceed to implement the above claim; that we can express the linear classifier in  $\phi$ -space without any geometric concepts except for the inner product. To do this, we insert equation 7 into equation 3 and 14, moving scalars out of the inner product brackets:

$$f(\mathbf{x}) = \text{sgn}\left(\sum_{i: \alpha_i \neq 0} (\alpha_i y_i \mathbf{x}_i^T \mathbf{x}) + b\right) \quad (15)$$

$$b = \frac{1}{|\{i : 0 < \alpha_i < C\}|} \sum_{i: 0 < \alpha_i < C} \left( y_i - \sum_{j: \alpha_j \neq 0} \alpha_j y_j \mathbf{x}_j^T \mathbf{x}_i \right) \quad (16)$$

Now, we are free to do the separation in any space, without explicit reference to that space. Just replace the inner product brackets in equations 4, 15 and 16 by the kernel function of your choice and the support vector machine will automatically carry out its job in the corresponding space.

## 4.2 Popular kernels

One common family of kernel functions is the polynomial kernels;  $K(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1^T \mathbf{x}_2 + 1)^d$ . Another is radial basis kernels;  $K(\mathbf{x}_1, \mathbf{x}_2) = e^{-\|\mathbf{x}_1 - \mathbf{x}_2\|^2 / (2\sigma^2)}$ . Linear separation in the original space is of course expressed by the kernel  $K(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^T \mathbf{x}_2$ . Which kernel to choose, and how to set the kernel parameters, depends on characteristics of the problem data.

### 4.3 Change your kernel

Rewrite your algorithm again, starting from the last version, the one that solved the dual problem relaxed to linearly non-separable data (section 3.1). Change the algorithm so that it uses the kernelised versions of equations 4, 15 and 16. Don't forget to substitute the kernel function  $K(\mathbf{x}_1, \mathbf{x}_2)$  for every inner product bracket  $\mathbf{x}_1^T \mathbf{x}_2$ . When you have found the  $\alpha_i$  you will no longer calculate an explicit  $\mathbf{w}$  vector; this vector now lives "hidden" in  $\phi$ -space. Instead, your classifier will be derived from equation 15. As you can see, the computation of this function requires passing the data points  $\mathbf{x}_i$  as well as the Lagrange multipliers  $\alpha_i$  and the bias  $b$ .

Use a polynomial kernel of low degree. You may pass a value for the degree parameter  $d$  to the kernel function as a global variable, if you like. (Hint: The implementation becomes easier if you pre-compute the matrix  $K$  as suggested in section 2.2.)

**Task:** Apply your algorithm to the linearly non-separable data. Find values for the parameters  $C$  and  $d$  that give you reasonable performance on the data set.

**Task:** Illustrate how your classifier divides up the plane using a contour plot, as before. Use equation 15 (without the  $\text{sgn}$ ) for the contour plot. Also plot the training examples, and mark the support vectors.

**Question (optional):** Tweak the  $C$  and  $d$  parameters. Report on how the two parameters affect classification performance, the number of support vectors ( $\alpha_i \neq 0$ ) and the way the plane is subdivided. For what values would you expect good generalisation performance?

*Good luck!*