

Lecture 6

*Lecturer: Amir Globerson**Scribe: Yishay Mansour*

6.1 Kernels for SVM classifiers

Today we continue with *Support Vector Machines (SVM)* and introduce the idea of a *kernel*. Using a kernel we will be able to find non-linear classifiers, but maintain computational efficiency.

We will cover the following topics:

- How many support vectors are there?
- Review of positive definite matrices
- Non linear classification and the kernel trick
- Characterization and construction of kernels.
- Optimization for SVM and stochastic gradient descent (details in next scribe).

6.1.1 Number of Support Vectors

When the data is linearly separable, we will expect at most $d + 1$ support vectors. There are degenerate cases where data exhibits very specific dependencies and more SVs are possible, but we generally do not expect this to happen. An example of a degenerate cases is when we only have two unique points (one per class) that are each repeated $\frac{m}{2}$ times. In this case all m points will be support vectors.

Note however, you might have less than $d + 1$ support vectors (See Figures 6.2 and 6.1 for an example on the plane).

To see why having $d + 1$ support vectors is sufficient, consider the SVM convex program,

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 \end{aligned}$$

The support vectors satisfy the constraint as an equality. This means that $[\mathbf{w}, b] \cdot [\mathbf{x}_i, 1] = 1$ for all support vectors. This is a linear equation in $d + 1$ dimensions. If we have more than $d + 1$ such equations we generally do not expect a solution (unless there is some degeneracy as above).

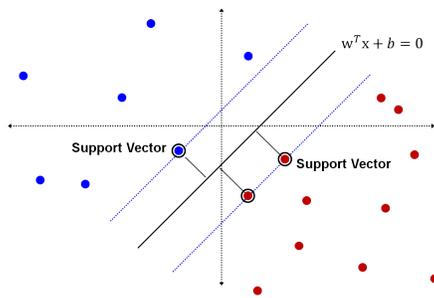


Figure 6.1: three support vectors

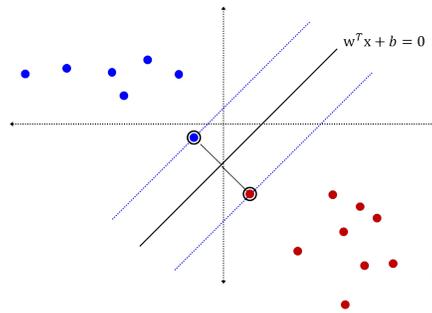


Figure 6.2: Two support vectors.

6.1.2 Positive Semi-definite (PSD) Matrices

A symmetric $d \times d$ matrix X is *positive semi-definite* (or PSD) if for any $\mathbf{v} \in \mathbb{R}^d$ such that $\mathbf{v} \neq \mathbf{0}$ we have $\mathbf{v}^T X \mathbf{v} \geq 0$. If strict inequality is required, the matrix is called positive definite (the $\mathbf{v} = 0$ case is excluded). We focus on PSD in what follows.

It can be shown (see slides) that a symmetric matrix X is positive definite if and only if there is a matrix B such that $X = BB^T$.

Given the above property, PSD matrices generalize the notion of a non-negative scalar, in the sense that they have a square root.

There are non-symmetric matrices, that can have the property that for any $\mathbf{v} \neq 0$ we have $\mathbf{v}^T X \mathbf{v} > 0$, for example $X = \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix}$. One can verify that $\mathbf{v}^T X \mathbf{v} = v_1^2 + v_2^2 > 0$. The test for a non-symmetric matrix X whether for any \mathbf{v} we have $\mathbf{v}^T X \mathbf{v} > 0$ is by testing if $X + X^T$ is PSD. Note that $X + X^T$ is symmetric. Also, for any \mathbf{v} , we have $\mathbf{v}^T X \mathbf{v} = \mathbf{v}^T X^T \mathbf{v}$, since it is a scalar. Therefore, $2\mathbf{v}^T X \mathbf{v} = \mathbf{v}^T (X + X^T) \mathbf{v}$, and X has the property iff $X + X^T$ is PSD.

It can also be shown that a matrix is PSD if and only if all its eigenvalues are non-negative (see slides).

6.1.3 The dual SVM formulation and PSD

Recall that the dual formulation of SVM is (when written as a minimization problem):

$$\begin{aligned}
 \min_{\alpha} \quad & \alpha^T \underbrace{\begin{bmatrix} y_1 y_1 \mathbf{x}_1 \cdot \mathbf{x}_1 & y_1 y_2 \mathbf{x}_1 \cdot \mathbf{x}_2 & \cdots & y_1 y_N \mathbf{x}_1 \cdot \mathbf{x}_N \\ y_2 y_1 \mathbf{x}_2 \cdot \mathbf{x}_1 & y_2 y_2 \mathbf{x}_2 \cdot \mathbf{x}_2 & \cdots & y_2 y_N \mathbf{x}_2 \cdot \mathbf{x}_N \\ \vdots & \vdots & \ddots & \vdots \\ y_N y_1 \mathbf{x}_N \cdot \mathbf{x}_1 & y_N y_2 \mathbf{x}_N \cdot \mathbf{x}_2 & \cdots & y_N y_N \mathbf{x}_N \cdot \mathbf{x}_N \end{bmatrix}}_M \alpha - \mathbf{1}^T \alpha \\
 \text{s.t.} \quad & \mathbf{y}^T \alpha = 0 \\
 & 0 \leq \alpha \leq C
 \end{aligned} \tag{6.1}$$

We can write the matrix M as PP^T where $P^T = [y_1 \mathbf{x}_1 \cdots y_m \mathbf{x}_m]$. For this reason the matrix M is positive semi-definite. This guarantees that the optimization is a convex quadratic program and can therefore be globally optimized efficiently.

6.1.4 Non-linear SVM

The basic SVM that we presented has a linear decision boundary. For example, in 1-dimension ($d = 1$) the decision boundary is simply a threshold. There are many cases when such a classifier isn't good enough. For example, consider Figure 6.3. Any threshold that we would select would have a high error rate. For this reason, introducing slack variables and using hinge loss, would not solve the problem. Any classifier we'll end up choosing will be using a threshold on the line, and any threshold is bad. An interesting solution is to map the points to a higher dimension, in this case 2-D. We can map each point x to (x, x^2) , namely map the line to a parabola. In figure 6.4 we show this for our example. Now there is a linear separator that can separate the data.

In general, we can always add enough dimensions such that the points are linearly separable, but in some extreme cases this might require a number of dimensions which is proportional to the number of points. In such a case the linear separator is very likely to overfit, unless it is done carefully as in the case of a RBF kernel. In terms of the generalization bounds we discussed, we are saying that an arbitrary mapping into n dimensional space will separate, but with a small margin, and therefore bad generalization error.

Another example is in figures 6.6 and 6.5. In 2D there is a circle which perfectly classifies the points. Namely, for each (x_1, x_2) the decision is based on $x_1^2 + x_2^2 \leq C$. We can get around this problem in a few ways. One solution is to move to polar coordinates, namely, map (x_1, x_2) to (r, θ) , where $r = \sqrt{x_1^2 + x_2^2}$ and $\tan(\theta) = x_2/x_1$.



Figure 6.3: An interval of blue between two red areas

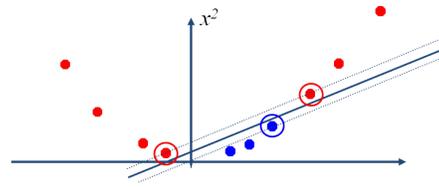


Figure 6.4: Lifting the interval to 2-D.

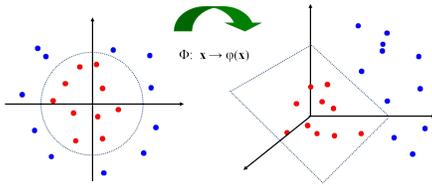


Figure 6.5: mapping from a circle to 3D

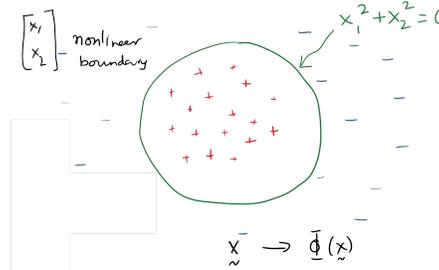


Figure 6.6: circle in 2-D.

An alternative solution, is to map (x_1, x_2) to (x_1^2, x_2^2) . This will solve our problem, but seems like a rather ad hoc solution. A better solution, which could handle any quadratic mapping, is to map (x_1, x_2) to $(1, x_1, x_2, x_1x_2, x_1^2, x_2^2)$ and then we can implement any quadratic function.

We can now discuss a general mapping of \mathbf{x} to $\phi(\mathbf{x})$. Let us consider the basic primal optimization.

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1 \quad i = 1, \dots, m \end{aligned}$$

The first question is *what is the dimension of the weights w ?* Note that this dimension depends on the dimension of $\phi(\mathbf{x})$. Since we would like to handle arbitrary large dimensions, this might be a computational problem. However, the rest of the math remains the same. If we have a very high dimensional feature vector $\phi(\mathbf{x})$ then the computational problem is twofold. First the inner products become an expensive operation, and second we need to compute and store a very large weight vector.

Let us consider the dual program in Equation 6.1. The good news here is that the new features $\phi(\mathbf{x})$ do not influence the constraint, and in the matrix M they appear only when computing an inner product of two generated by $\phi(\mathbf{x})$. Let us consider the hypothesis h

that we build. Recall that $\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \phi(\mathbf{x}_i)$ and then

$$h(\mathbf{x}) = \mathbf{w} \cdot \phi(\mathbf{x}) + b = \sum_{i=1}^m \alpha_i y_i \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}) + b$$

again, we need to consider only inner products of $\phi(\mathbf{x})$.

Finally, we compute b by considering any support vector (\mathbf{x}, y) and setting

$$b = y - \sum_{i=1}^m \alpha_i y_i \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x})$$

again, we need to consider only inner products of $\phi(\mathbf{x})$.

This suggests that if we can compute the inner product of feature vectors $\phi(\mathbf{x})$ efficiently, we can perform the entire mapping efficiently. This is precisely the idea behind the *kernel trick*. Rather than computing the inner product explicitly, we can compute it implicitly. For a mapping $\phi(\mathbf{x})$ we define a kernel K such that

$$K(\mathbf{x}', \mathbf{x}'') = \phi(\mathbf{x}')^T \phi(\mathbf{x}'')$$

and then the hypothesis becomes

$$h(\mathbf{x}) = \sum_{i=1}^m \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$

Also in the matrix M we can replace any $\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$ by $K(\mathbf{x}_i, \mathbf{x}_j)$.

Going back to the example of figure 6.4. We can set $\phi(\mathbf{x}) = (x, x^2)$ and then

$$K(x, z) = \phi(\mathbf{x}) \cdot \phi(\mathbf{z}) = xz + x^2 z^2$$

In this example we did not gain much by not doing the inner product explicitly using ϕ . However, in higher dimensions there will be a significant benefit.

For the quadratic polynomial kernel, we can first describe the kernel itself, and only then show that indeed it is an inner product of some ϕ . The quadratic polynomial K_2 is,

$$K_2(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x} \cdot \mathbf{x}')^2$$

clearly this kernel can be computed in time $O(d)$ for $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$. We still need to show that there are functions ϕ such that $K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x}) \cdot \phi(\mathbf{x}')$. We can compute the kernel

explicitly. Indeed for $d = 2$ we have:

$$\begin{aligned} K_2(\mathbf{x}, \mathbf{x}') &= (1 + \mathbf{x} \cdot \mathbf{x}')^2 = (1 + x_1x'_1 + x_2x'_2)^2 \\ &= 1 + x_1^2(x'_1)^2 + x_2^2(x'_2)^2 + 2x_1x_2x'_1x'_2 + 2x_1x'_1 + 2x_2x'_2 \\ &= \begin{pmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ \sqrt{2}x'_1 \\ \sqrt{2}x'_2 \\ x_1'^2 \\ x_2'^2 \\ \sqrt{2}x'_1x'_2 \end{pmatrix} \\ &= \phi(\mathbf{x}) \cdot \phi(\mathbf{x}') \end{aligned}$$

In a similar way this can be extended to $\mathbf{x} \in \mathbb{R}^d$. The main benefit is that the computation is $O(d)$ while the length of ϕ is $O(d^2)$.

There are several classes of popular kernels. The first is the *linear kernel* which is the basic SVM. We have the degree r polynomial kernel K_r , where,

$$K_r(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x} \cdot \mathbf{x}')^r$$

The degree r kernel uses a feature vector of length $O(d^r)$ and the computation is only $O(d)$. The radial basis function (RBF) kernel K_g has a parameter σ and

$$K_g(\mathbf{x}, \mathbf{x}') = e^{-\|\mathbf{x}-\mathbf{x}'\|^2/(2\sigma^2)}$$

The feature vector of the RBF kernel is of *infinite* size (see slides). Still the computation of a RBF kernel is $O(d)$.

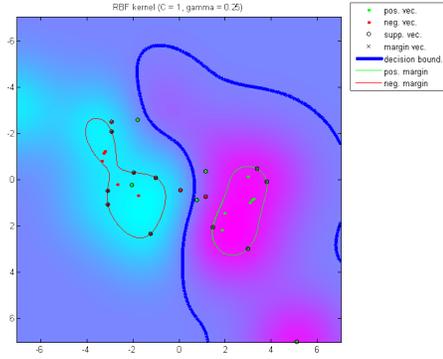
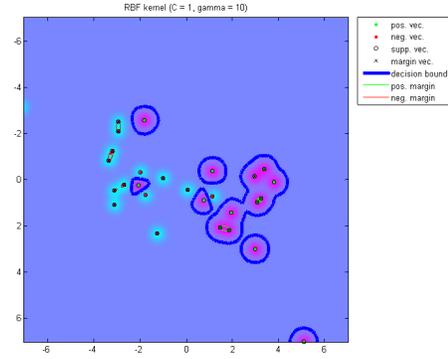
To better understand the role of a kernel we can define a kernel matrix K_S where $[K_S]_{i,j} = \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j) = K(\mathbf{x}_i, \mathbf{x}_j)$. Note that the matrix depends on the points $\mathbf{x}_1, \dots, \mathbf{x}_n$.

With the RBF kernel we can select the points \mathbf{x}_i such that $K(\mathbf{x}_i, \mathbf{x}_j) \leq \epsilon$ for $i \neq j$, and for sufficiently small ϵ , the matrix K_S is full rank. This implies that the feature vector grows with the number of points, and hence is infinite.

The RBF Radial Basis Function (RBF) kernels are very flexible. They are often parameterized by $\gamma = 1/(2\sigma)$ rather than σ . For large values of γ they are essentially a nearest-neighbor computation since the influence falls exponentially with the distance (see Figures 6.8 and 6.7).

6.1.5 Characterizing Kernels

How do we find a function $K(\mathbf{x}, \mathbf{x}')$ that corresponds to a dot product in *some* feature space $\phi(\mathbf{x})$? To formalize the result, we first define *Positive Definite Kernels*. These are functions $K(\mathbf{x}, \mathbf{x}')$ that satisfy:

Figure 6.7: RBF kernel with large σ Figure 6.8: RBF kernel with small σ .

1. Symmetry: $K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}', \mathbf{x})$.
2. For any sample $S = [\mathbf{x}_1, \dots, \mathbf{x}_m]$ the matrix K_S is PSD.

It can be shown that a kernel is a dot product in feature space if and only if it is Positive Definite Kernel. This result is known as Mercer's theorem (Mercer's result is typically stated for an equivalent notion of positive definiteness).

We can compose PD kernels from existing ones. Here are two simple examples:

1. Addition $K_3(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}') + K_2(\mathbf{x}, \mathbf{x}')$. If ϕ_1 and ϕ_2 are the feature vectors of K_1 and K_2 , we can create $\phi_3(\mathbf{x}) = [\phi_1(\mathbf{x}); \phi_2(\mathbf{x})]$, namely concatenating the two feature vectors.
2. Multiplication: $K_3(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}') \cdot K_2(\mathbf{x}, \mathbf{x}')$. If ϕ_1 and ϕ_2 are the feature vectors of K_1 and K_2 we can do the following:

$$\begin{aligned}
 K_3(\mathbf{x}, \mathbf{x}') &= (\phi_1(\mathbf{x})^T \phi_1(\mathbf{x}'))(\phi_2(\mathbf{x}')^T \phi_2(\mathbf{x})) \\
 &= \text{tr}(\phi_1(\mathbf{x})^T \phi_1(\mathbf{x}') \phi_2(\mathbf{x}')^T \phi_2(\mathbf{x})) \\
 &= \text{tr}(\phi_2(\mathbf{x}) \phi_1(\mathbf{x})^T \phi_1(\mathbf{x}') \phi_2(\mathbf{x}')^T) \\
 &= \langle \text{VEC}(\phi_2(\mathbf{x}) \phi_1(\mathbf{x})^T), \text{VEC}(\phi_1(\mathbf{x}') \phi_2(\mathbf{x}')^T) \rangle,
 \end{aligned}$$

where we used the fact that $\text{tr}(ABCD) = \text{tr}(DABC)$ and $\text{VEC}(\cdot)$ transforms a matrix to a vector.

Suppose we want to do Nearest Neighbor in feature space, how can we use kernels. Recall that

$$\|a - b\|^2 = (a - b)^T(a - b) = a^T a - 2a^T b + b^T b$$

Similarly, in feature space we have

$$\|\phi(a) - \phi(b)\|^2 = K(a, a) - 2K(a, b) + K(b, b)$$

We can create a kernel from almost any proper distance metric. One way of doing it is using the *generalized RBF Kernels*, where

$$K(h_1, h_2) = e^{-D^2(h_1, h_2)/\beta}$$

There are a few ways of setting the distance $D(h_1, h_2)$:

- L_1 distance: $D(h_1, h_2) = \sum_{i=1}^d |h_1(i) - h_2(i)|$
- L_2 distance: $D^2(h_1, h_2) = \sum_{i=1}^d (h_1(i) - h_2(i))^2$
- L_∞ distance: $D(h_1, h_2) = \max_{i=1}^d |h_1(i) - h_2(i)|$
- χ^2 distance: $D^2(h_1, h_2) = \sum_{i=1}^d \frac{(h_1(i) - h_2(i))^2}{h_1(i) + h_2(i)}$ which is very applicable to histograms.
- Hellinger distance: $D^2(h_1, h_2) = \sum_{i=1}^d (\sqrt{h_1(i)} - \sqrt{h_2(i)})^2$ which is very applicable to distributions.
- Mahalanobis distance: $D^2(h_1, h_2) = (h_1 - h_2)^T S^{-1} (h_1 - h_2)$ which allows to learn the metric, which is parametrized by a positive definite matrix S .