

CSE 252C: Computer Vision III

Lecturer: Serge Belongie

Scribe: Catherine Wah

LECTURE 15

Kernel Machines

15.1. Kernels

We will study two methods based on a special kind of function $k(\mathbf{x}, \mathbf{y})$ called a *kernel*:

- Kernel PCA (to help build our intuition) and
- Support Vector Machines, which are a powerful discriminative learning method.

What is a kernel? In this context, it is a function that satisfies the following condition:

$$(15.1) \quad \sum_{i,j} k(\mathbf{x}^i, \mathbf{x}^j) c_i c_j \geq 0,$$

for some set of feature vectors \mathbf{x} and any vector \mathbf{c} . We say the kernel is positive semi-definite, or that it satisfies Mercer's condition.

¹Department of Computer Science and Engineering, University of California, San Diego.

15.2. Kernel Trick

What’s interesting about this, and what forms the basis of the so-called “kernel trick” (Aizerman *et al.*, 1964), is that kernels that satisfy this condition are equivalent to a simple dot product in some higher dimensional space:

$$(15.2) \quad k(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{y}),$$

where $\phi(\cdot)$ maps a vector from “input space” to “feature space” (usually higher, or infinite, dimensional). In some cases we specify $\phi(\cdot)$, in other cases, it’s implicit for a choice of $k(\cdot, \cdot)$. The “kernel trick” refers to the idea of swapping in $k(\mathbf{x}, \mathbf{y})$ for $\mathbf{x} \cdot \mathbf{y}$ in algorithms that call for a dot product between two vectors, for example, the perceptron, SVMs (supervised), PCA (unsupervised).

In some cases, it may not be obvious that the algorithms, at their heart, come down to dot products, but these and several others have been formulated in such a way to exhibit this characteristic. So what does this buy us? In the supervised case, classes that aren’t linearly separable in input space might be linearly separable in feature space (Figure 1). Consider the

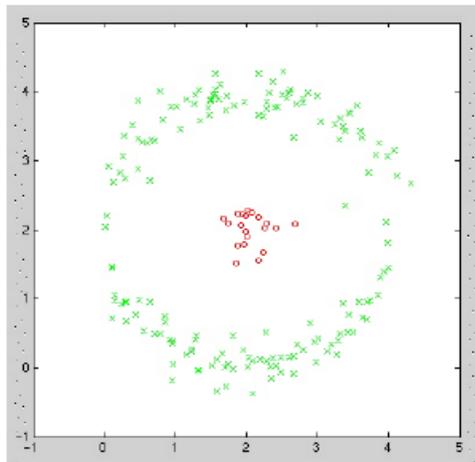


Figure 1. The clump and annulus problem. These classes are linearly inseparable in the input space. Note: a quadratic surface would do the trick here, but linear algorithms are simpler and better understood.

mapping $\Phi : (x_1, x_2)^T \mapsto (x_1, x_2, x_1^2 + x_2^2)^T$, where $\mathbf{x} \in \mathbb{R}^2$ and $\phi(\mathbf{x}) \in \mathbb{R}^3$. In this way, by adding a new dimension, it becomes possible to separate the clump from the annulus.

This hints to us that high dimensionality can be good (greater classification power)! Your instinct might be telling you that high dimensionality can also get you into trouble, conceptually or computationally; this instinct is correct.

This idea is the subject of Vapnik-Chervonenkis (VC) theory, a critical component of statistical learning theory, which addresses questions such as how best to perform classification given nothing but labelled example data (with no prior knowledge). We won't study VC theory in this class, but it provides us with answers to questions such as "under what conditions are high dimensional representations good?"

Supposing we restrict ourselves to cases in which the high dimensional mapping is good, there is still the problem of computation time. For example, a simple polynomial mapping on 28×28 MNIST digits can easily explode the dimensionality into the billions. This is where the kernel trick comes in. Except for these toy examples, in practice, one never makes explicit use of $\Phi(\cdot)$, or even necessarily needs to know what it is. We only need to know that it exists, which Mercer's condition answers for us.

This was foreshadowed in the second homework, where we considered the Mahalanobis distance between two vectors:

$$(15.3) \quad (\mathbf{x}^i - \mathbf{x}^j)^\top \Sigma^{-1} (\mathbf{x}^i - \mathbf{x}^j) = (\mathbf{y}^i - \mathbf{y}^j)^\top (\mathbf{y}^i - \mathbf{y}^j) = \|\mathbf{y}^i - \mathbf{y}^j\|^2$$

where, if Σ^{-1} is positive semi-definite, we can write $\Sigma^{-1} = QQ^\top$ and set $\mathbf{y} = Q^\top \mathbf{x}$. In this case, Q represents a simple case of a feature mapping; the Mahalanobis distance is equivalent to a dot product in a different space.

15.3. Kernel PCA

Now let's look at an application of the kernel trick: Kernel PCA (Schölkopf *et al.*, 1998). This is an unsupervised example, used for dimensionality reduction when the data doesn't live on a linear manifold. We use the kernel trick to perform "classic" PCA in high dimensional space, then map it back to input space, where the "axes" will appear curved (Figure 2).

Before we "kernelize" PCA, let's review how classic PCA works (assuming centered data: $\sum_{i=1}^N \mathbf{x}^i = \mathbf{0}$). PCA finds the eigenvectors of the covariance matrix C :

$$(15.4) \quad C = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}^i)(\mathbf{x}^i)^\top.$$

By construction, C is positive semidefinite. Eigenvectors $\lambda \mathbf{v} = C \mathbf{v}$ are used to capture principal axes of variation.

We can use C in the following way:

- (a) Find the eigenvectors and eigenvalues, sorted in decreasing order by eigenvalue, possibly truncating
- (b) Project test data onto eigenvectors
- (c) Use the projections for classification, denoising, compression, etc.

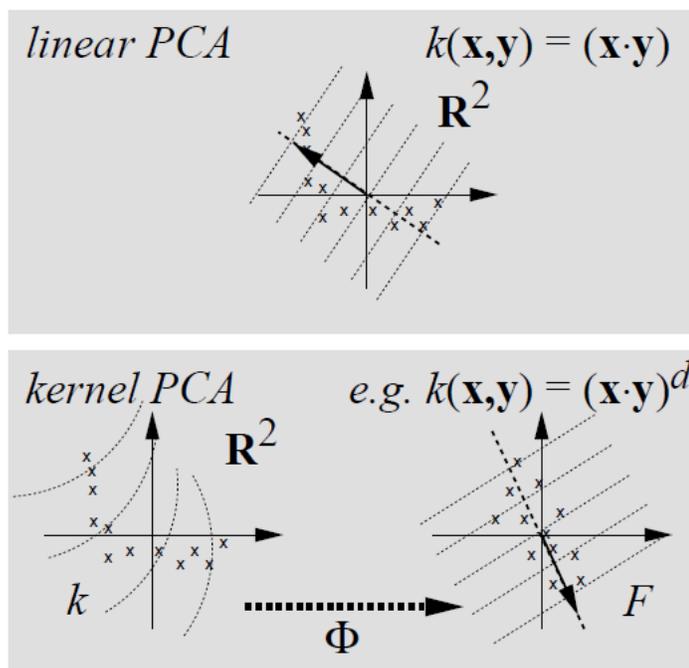


Figure 2. The basic idea of kernel PCA with 2D inputs. In some high-dimensional feature space F (bottom right), we are performing linear PCA, just like a PCA in input space (top). Since F is nonlinearly related to the input space (via Φ), the contour lines of constant projections onto the principal eigenvector become nonlinear in input space (Schölkopf *et al.*, 1998).

We'd like to be able to do this for nonlinear data, therefore, the solution is to kernelize. In order to do so, we start by expressing PCA in terms of dot products:

$$(15.5) \quad C\mathbf{v} = \frac{1}{N} \sum_i \mathbf{x}^i \mathbf{x}^{i\top} \mathbf{v} = \lambda \mathbf{v},$$

therefore

$$(15.6) \quad \mathbf{v} = \frac{1}{N\lambda} \sum_j \mathbf{x}^j \mathbf{x}^{j\top} \mathbf{v} = \frac{1}{N\lambda} \sum_j (\mathbf{x}^j \cdot \mathbf{v}) \mathbf{x}^j.$$

(Note here that $(\mathbf{x}^j \cdot \mathbf{v})$ is just a scalar.) By using dot products, all solutions \mathbf{v} with $\lambda \neq 0$ lie in the span of $\{\mathbf{x}^i\}$, or $\mathbf{v} = \sum_i \alpha_i \mathbf{x}^i$, where the α_i 's are the dot products.

Now suppose we pick some mapping $\phi(\cdot)$. Assume we center the data (a homework problem), and as before, we write

$$(15.7) \quad C = \frac{1}{N} \sum_i \phi(\mathbf{x}^i) \phi(\mathbf{x}^i)^\top,$$

which is possibly high dimensional, and we diagonalize it as $\lambda \mathbf{v} = C \mathbf{v}$. Like before, \mathbf{v} lies in the span of $\phi(\mathbf{x}^i)$'s.

Now we dot both sides by $\phi(\mathbf{x}^j)$:

$$(15.8) \quad \lambda(\phi(\mathbf{x}^j) \cdot \mathbf{v}) = \phi(\mathbf{x}^j) \cdot C \mathbf{v}, \forall j$$

and note that $\mathbf{v} = \sum_i \alpha_i \phi(\mathbf{x}^i)$. The motivation for this is that we want to avoid explicit use of $\phi(\cdot)$; we only want to use it via $k(\cdot, \cdot)$.

Therefore,

$$(15.9) \quad \lambda \sum_i \alpha_i (\phi(\mathbf{x}^j) \cdot \phi(\mathbf{x}^i)) = \frac{1}{N} \sum_i \alpha_i \left(\phi(\mathbf{x}^j) \cdot \sum_l \phi(\mathbf{x}^l) \right) (\phi(\mathbf{x}^l) \cdot \phi(\mathbf{x}^i)), \forall k$$

or

$$(15.10) \quad N \lambda K \boldsymbol{\alpha} = K^2 \boldsymbol{\alpha}$$

where $K_{ij} = \phi(\mathbf{x}^i) \cdot \phi(\mathbf{x}^j)$ and $\boldsymbol{\alpha}$ is a vector of α_i 's. Note that since K is full rank, it follows that

$$(15.11) \quad N \lambda \boldsymbol{\alpha} = K \boldsymbol{\alpha}.$$

This is our eigenvector problem, now on K instead of the covariance matrix.

So in the end, kernel PCA just requires us to diagonalize K instead of C , a different (and bigger) covariance matrix defined on vectors in a high dimensional space. In practice, the data needs to be centered first, which can be done via a simple operation on K ; see Homework 4.

How do we use these principal components? Consider some test point \mathbf{x} (*e.g.*, any point in \mathbb{R}^2). To project it onto eigenvectors, we compute:

$$(15.12) \quad \mathbf{v}^k \cdot \phi(\mathbf{x}) = \sum_i \alpha_i k(\mathbf{x}^i, \mathbf{x})$$

for the k th kernel principal component coefficient. We demonstrate this for a 2D toy example with 3 gaussian clumps, using a gaussian kernel (Figure 3). Recall that regular PCA would just give two orthogonal axes in \mathbb{R}^2 . Surface brightness shows the value of the k th eigenvector projection; for linear PCA you would just have 2 ramps. Note the similarity to spectral clustering; NCut used $D^{-1/2} W D^{-1/2}$ instead of centering in feature space. It allows extrapolation to the full plane; this is the same thing as the Nyström extension, which is used in spectral clustering. Note in this example, each cluster is first “found,” then a local x - y coordinate system is extracted for

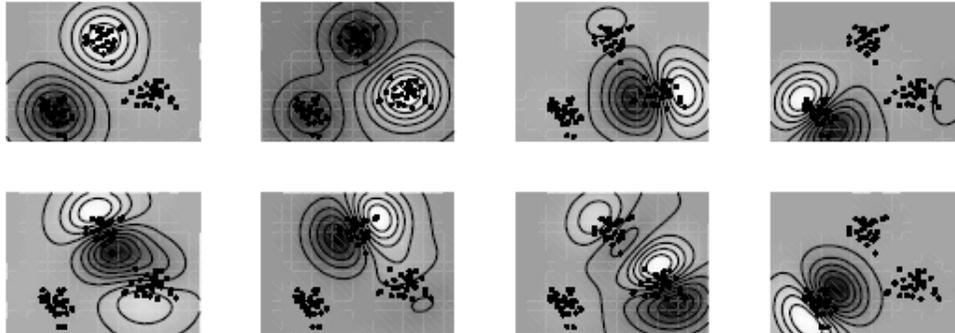


Figure 3. Toy example with three data clusters; first eight nonlinear principal components are extracted with a radial basis function (*i.e.*, Gaussian) for the kernel. Note that the first two principal components (top left) nicely separate the three clusters (Schölkopf *et al.*, 1998).

it. For denoising, we use a fixed point iteration to find the preimage; this is related to mean-shift.

15.4. Support Vector Machines

Now let's look at a discriminative example: the support vector machine (SVM). Now we have data labeled (positive and negative):

$$(15.13) \quad \{(\mathbf{x}^1, y^1), (\mathbf{x}^2, y^2), \dots, (\mathbf{x}^N, y^N)\}, y \in \{+1, -1\}$$

In the linear case, we want to find a separating hyperplane with maximum margin, *i.e.*, we want a hyperplane with maximal distance to the nearest data point (Figure 4). The linear classifier (decision boundary) has the form:

$$(15.14) \quad \mathbf{w} \cdot \mathbf{x} - b = 0,$$

where \mathbf{w} is the normal vector for the hyperplane and b is the offset.

SVM learning requires the solution of a quadratic programming problem (beyond the scope of this class) that returns the data points that “prop up” the hyperplanes that define the margin, parallel to the separating hyperplane. This can be generalized to the soft-margin case, when perfect separation is not possible.

The extension to cases requiring a nonlinear separation boundary requires kernelizing; the linear separating boundary looks curved in input space. When testing what side of the decision boundary you are on, we use a kernel in place of the dot product, which we can think of as the distance (or similarity) to exemplars right on either side of the margin; *e.g.*, for face gender classification, these are highly hermaphroditic faces $\pm\epsilon$.

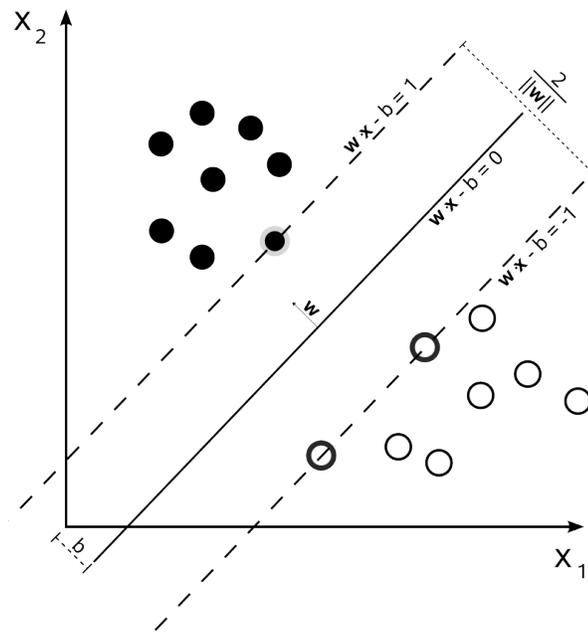


Figure 4. Maximum-margin hyperplane and margins for a SVM trained with samples from two classes. Samples on the margin are called the support vectors. (http://en.wikipedia.org/wiki/Support_vector_machine)

To summarize, SVMs are a highly general method with excellent results. The usual questions of what kernel to use, what value of σ , etc. can be addressed with experience and cross validation.