

# Advanced Machine Learning Subsidiary Notes

## Lecture 14: Kernel Trick

Adam Prügel-Bennett

April 23, 2020

## 1 Keywords

- The Kernel Trick, SVMs, Regression

## 2 Main Points

### 2.1 Kernels

- SVM kernels are symmetric functions of two variable that can be factorised as an inner-product

$$K(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^\top \phi(\mathbf{y}) = \sum_{i=1}^{p'} \phi_i(\mathbf{x}) \phi_i(\mathbf{y})$$

- $\phi(\mathbf{x})$  are vectors whose elements,  $\phi_i(\mathbf{x})$ , are real-valued functions of the features  $\mathbf{x}$  (every different feature will correspond to a different vector  $\phi(\mathbf{x})$ )
- $p'$  is the dimensionality of the extended feature space which might be infinite
- An immediate consequence of this is that the vectors are positive semi-definite
  - \* This follows because for any function  $f(\mathbf{x})$  the quadratic form is non-negative

$$\begin{aligned} \iint f(\mathbf{x}) K(\mathbf{x}, \mathbf{y}) f(\mathbf{y}) d\mathbf{x} d\mathbf{y} &= \iint f(\mathbf{x}) \phi(\mathbf{x})^\top \phi(\mathbf{y}) f(\mathbf{y}) d\mathbf{x} d\mathbf{y} \\ &= \sum_{i=1}^{p'} \left( \int f(\mathbf{x}) \phi_i(\mathbf{x}) d\mathbf{x} \right)^2 \geq 0 \end{aligned}$$

- **Eigenfunctions and Mercer's Theorem**

- Kernel functions play the same role for functions as matrices do for normal vectors
  - \* that is they describe general linear transformations
  - \* for a function  $f(\mathbf{x})$  the argument  $\mathbf{x}$  can be seen as an index just like  $i$  is the index of element  $v_i$  of a vector  $\mathbf{v}$
  - \* we will consider only symmetric kernels (that is, where  $K(\mathbf{x}, \mathbf{y}) = K(\mathbf{y}, \mathbf{x})$ )
  - \* these play a similar role as symmetric matrices
- Eigensystems for Kernels
  - \*  $\psi(\mathbf{y})$  is said to be an eigenfunction of a kernel functions if

$$\int K(\mathbf{x}, \mathbf{y}) \psi(\mathbf{y}) d\mathbf{y} = \lambda \psi(\mathbf{x})$$

- \* In an analogy to the eigen-decomposition of a symmetric matrix we can define the eigen-decomposition of a symmetric kernel function

$$K(\mathbf{x}, \mathbf{y}) = \sum_i \lambda_i \psi_i(\mathbf{x}) \psi_i(\mathbf{y})$$

- \* This is known as **Mercer's Theorem**
- \* We proved the decomposition for matrices
  - the difficult part in the proof is that you need the eigenvectors to span the vector space
  - this is intuitively obvious if there are  $n$  orthogonal eigenvectors in an  $n$ -dimensional space
  - it is harder in functions spaces and you need to define the vector space you are modelling (e.g.  $L_2$ )
  - if you assume that the set of eigenvectors span the function space then the rest of the proof is the same as for matrices
  - don't worry if you don't understand this it is enough to remember Mercer's Theorem
- \* Mercer's Theorem holds for any symmetric kernel function (it does not have to be positive semi-definite)
- \* But if  $K(\mathbf{x}, \mathbf{y})$  are positive semi-definite then there exist real functions  $\phi_i(\mathbf{x}) = \sqrt{\lambda_i} \psi_i(\mathbf{x})$  such that

$$K(\mathbf{x}, \mathbf{y}) = \sum_i \lambda_i \psi_i(\mathbf{x}) \psi_i(\mathbf{y}) = \sum_i \phi_i(\mathbf{x}) \phi_i(\mathbf{y})$$

- if  $K(\mathbf{x}, \mathbf{y})$  was not positive semi-definite then some of the eigenvalues would be negative and the functions  $\psi_i(\mathbf{x})$  would not be real-valued

## 2.2 SVM Kernels

- SVM Kernels are positive semi-definite symmetric functions
  - There are four necessary and sufficient conditions that hold for any positive semi-definite kernel
    1. All their eigenvalues are non-negative (i.e. either zero or positive)
    2. They can be decomposed as

$$K(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^\top \phi(\mathbf{y}) = \sum_{i=1}^{p'} \phi_i(\mathbf{x}) \phi_i(\mathbf{y})$$

where  $\phi_i(\mathbf{x})$  are real-valued functions

3. Their quadratic form with any function  $f(\mathbf{x})$  is non-negative
4. For any set of points  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  the matrix

$$\mathbf{K} = \begin{pmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & K(\mathbf{x}_1, \mathbf{x}_2) & \cdots & K(\mathbf{x}_1, \mathbf{x}_n) \\ K(\mathbf{x}_2, \mathbf{x}_1) & K(\mathbf{x}_2, \mathbf{x}_2) & \cdots & K(\mathbf{x}_2, \mathbf{x}_n) \\ \vdots & \vdots & \ddots & \vdots \\ K(\mathbf{x}_n, \mathbf{x}_1) & K(\mathbf{x}_n, \mathbf{x}_2) & \cdots & K(\mathbf{x}_n, \mathbf{x}_n) \end{pmatrix}$$

is a positive semi-definite matrix

- \* such matrices are known as **Gram matrices**
- \* I didn't mention this in the lecture and won't use this property, but for completeness I mention it here (you won't be tested on it)
- \* the proof that this is a necessary condition follows rather simply from the fact that if we define a matrix  $\Phi$  with elements  $\Phi_{ik} = \phi_i(\mathbf{x}_k)$  then  $\mathbf{K} = \Phi^\top \Phi$  and we have seen many times any such matrix is positive semi-definite
- Recall from the previous lecture that any kernel function that allows a decomposition in terms of positive functions can be used an SVM where we can use the kernel trick
  - If we don't use positive semi-definite kernels then our "distances" (used in computing margins) are no-longer proper distances and can be negative (invalidating everything)

## 2.3 Constructing SVM Kernel

- Most functions of two variable won't be positive semi-definite
- Given a function of two variables it is hard to determine if it is positive-semi definite (none of the definitions are particularly easy to use)
- However we can use simple rules to build positive-semi definite (PSD) kernels from other positive semi-definite kernels

1. Our starting point is to note the inner produce  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y}$  is positive semi-definite
  - as an aside we don't necessarily need to use normal vectors as our features so long as we objects with an inner-product

2. Adding PSD kernels

*if  $K_1(\mathbf{x}, \mathbf{y})$  and  $K_2(\mathbf{x}, \mathbf{y})$  are PSD kernels then so is  $K_3(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x}, \mathbf{y}) + K_2(\mathbf{x}, \mathbf{y})$*

- To prove this we can use the property that PSD have non-negative quadratic form

$$\begin{aligned} Q &= \int f(\mathbf{x}) K_3(\mathbf{x}, \mathbf{y}) f(\mathbf{y}) d\mathbf{x} d\mathbf{y} \\ &= \int f(\mathbf{x}) (K_1(\mathbf{x}, \mathbf{y}) + K_2(\mathbf{x}, \mathbf{y})) f(\mathbf{y}) d\mathbf{x} d\mathbf{y} \\ &= \int f(\mathbf{x}) K_1(\mathbf{x}, \mathbf{y}) f(\mathbf{y}) d\mathbf{x} d\mathbf{y} + \int f(\mathbf{x}) K_2(\mathbf{x}, \mathbf{y}) f(\mathbf{y}) d\mathbf{x} d\mathbf{y} \geq 0 \end{aligned}$$

3. Multiplication by a positive scalar

*if  $K_1(\mathbf{x}, \mathbf{y})$  is a PSD kernels and  $c > 0$  then so is  $K_3(\mathbf{x}, \mathbf{y}) = c K_1(\mathbf{x}, \mathbf{y})$*

- We can prove this in a similar way to the last proof

4. Multiply PSD kernels

*if  $K_1(\mathbf{x}, \mathbf{y})$  and  $K_2(\mathbf{x}, \mathbf{y})$  are PSD kernels then so is  $K_3(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x}, \mathbf{y}) K_2(\mathbf{x}, \mathbf{y})$*

- This is easy to prove using the decomposition of PSD to inner products

$$K_3(\mathbf{x}, \mathbf{y}) = \sum_{i,j} \phi_i^1(\mathbf{x}) \phi_i^1(\mathbf{y}) \phi_j^2(\mathbf{x}) \phi_j^2(\mathbf{y}) = \sum_{i,j} \phi_{ij}^3(\mathbf{x}) \phi_{ij}^3(\mathbf{y})$$

where  $\phi_{ij}^3(\mathbf{x}) = \phi_i^1(\mathbf{x}) \phi_j^2(\mathbf{x})$

\* this (double) sum we can treat as an inner-product

\* it is easy to show that the quadratic form with any function  $f(\mathbf{x})$  is non-negative

5. Powers of PSD kernels

*if  $K_1(\mathbf{x}, \mathbf{y})$  is a PSD kernels then so is  $K_1^n(\mathbf{x}, \mathbf{y})$  for any natural number  $n$*

- Since the product of any two PSD kernels are PSD then the square of a PSD kernel is PSD
- But by an inductive argument this holds for any integer power

6. Exponential of PSD kernels

*The exponential of a PSD kernel is also a PSD kernel*

- convergent Taylor expansions allow us to approximate a function to any degree of accuracy
- often Taylor expansions aren't everywhere convergent (so we have to be careful)
- but Taylor expansions of exponentials are everywhere convergent
- further Taylor expanding an exponential of a PSD kernel involves a sum of PSD kernels

$$e^{K(\mathbf{x}, \mathbf{y})} = \sum_i \frac{1}{i!} K^i(\mathbf{x}, \mathbf{y}) = 1 + K(\mathbf{x}, \mathbf{y}) + \frac{1}{2} K^2(\mathbf{x}, \mathbf{y}) + \dots$$

\* each term is a PSD kernel

- Using these properties we see that  $K(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|^2}$  is a PSD kernel if  $\gamma > 0$ 
  - Since
$$e^{-\gamma \|\mathbf{x} - \mathbf{y}\|^2} = e^{-\gamma \|\mathbf{x}\|^2} e^{2\gamma \mathbf{x}^\top \mathbf{y}} e^{-\gamma \|\mathbf{y}\|^2}$$
  - But  $e^{-\gamma \|\mathbf{x}\|^2}$  and  $e^{-\gamma \|\mathbf{y}\|^2}$  are just positive constants
  - $\mathbf{x}^\top \mathbf{y}$  is an inner product so a PSD kernel
  - Since  $2\gamma > 0$  then  $2\gamma \mathbf{x}^\top \mathbf{y}$  is a PSD kernel
  - But then so is  $e^{2\gamma \mathbf{x}^\top \mathbf{y}}$
  - This kernel is known as the *radial basis function* or *RBF* or *Gaussian kernel*
  - It has a hyper-parameter,  $\gamma$  that determines the length scale in the problem (or rather inverse-length scale)
  - this is a very important kernel as it often (but certainly not always) gives good performance (if  $\gamma$  is appropriately chosen)
- **Non-numerical Kernels**
  - When SVMs were fashionable there was a whole industry of researchers finding clever kernels
  - When working with language or trees or graphs it paid to create bespoke kernels for these structures
  - Typically these would all be built up from inner-products
  - Using clever algorithms you can build very clever kernels functions
  - One down side of SVM kernels is they don't naturally capture prior knowledge about the problem being tackled
    - \* a clever work around is to build SVMs based on other learning machines that are trained the problem
    - \* an example of this is the use of *Fisher kernels* based on Fisher information

## 2.4 Beyond SVMs

- There are a lot of other kernel based learning machines
- Many of these use constraints
- They often involve linear operations between vectors where the optimum depends on the inner-product of vectors
  - thus we can use the kernel trick
- *SVR* are support vector machines for regression
  - here we try to find a dividing plane so that all points lie within a margin (the exact opposite of what we had)
  - We can introduce slack variables to allow some points to lie outside the margin
    - \* the slack variables must be non-negative
    - \* we can use a linear punishment  $s_i$  or quadratic punishment  $s_i^2$
- We can also do *kernel ridge regression*

$$\min_{\mathbf{w}} \lambda \|\mathbf{w}\|^2 + \sum_i \left( y_i - \mathbf{w}^\top \phi(\mathbf{x}_i) \right)^2$$

- $\|\mathbf{w}\|^2$  is a regularisation term

- The weights must lie in the space spanned by the set of extended feature vectors  $\{\phi(\mathbf{x}_k) | k = 1, 2, \dots, m\}$
- Thus we can write

$$\mathbf{w} = \sum_i \alpha_i \phi(\mathbf{x}_i)$$

- \* Note that here  $\alpha_i$  are just parameters; they are not Lagrange multipliers and they can be negative
- Substituting this into the objective function for ridge regression we get a quadratic optimisation problem in  $\boldsymbol{\alpha}$  that just depends on the inner products  $\phi^\top(\mathbf{x}_i) \phi(\mathbf{x}_j)$
- We can use the kernel trick
- *Kernel PCA*
  - For kernel PCA we map features into an external feature space
  - We then use the dual form of PCA (which we've done in an earlier lecture)
  - This allows us to find non-linear manifolds where the data varies
- *Kernel Canonical Correlation Analysis*
  - Canonical correlation analysis finds correlations between datasets
  - The linear form is a bit naff
  - But the kernel form can give nice results
- *Gaussian Processes*
  - Gaussian Processes also use kernels
  - They are a bit different to other kernel methods
    - \* we don't think of the working in an extended feature space
    - \* but they are PSD
  - They are one of the most successful methods for doing regression
  - We will look at them later

## 3 Exercises

### 3.1 Quadratic Kernels

- Show that the kernel function  $K(\mathbf{x}, \mathbf{y}) = \phi^\top(\mathbf{x}) \phi(\mathbf{y})$ , where

$$\phi(\mathbf{x}) = (x_1^2, x_2^2, x_3^2, \sqrt{2}x_1x_2, \sqrt{2}x_1x_3, \sqrt{2}x_2x_3)$$

can be written as  $(\mathbf{x}^\top \mathbf{y})^2$  if  $\mathbf{x}$  and  $\mathbf{y}$  are vectors of length 3.

- Answer below

### 3.2 Kernel Ridge Regression

- Work out the details for kernel ridge regression
- Have a go at implementing kernel ridge regression on a real data set
- I'll leave you to work this out

## 4 Experiments

### 4.1 Gram Matrix

- Generate ten random vectors  $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{10})$  where  $\mathbf{x}_k \in \mathbb{R}^5$
- Compute the Gram matrix  $\mathbf{K}$  with components

$$K_{kl} = K(\mathbf{x}_k, \mathbf{x}_l) = e^{-\|\mathbf{x}_k - \mathbf{x}_l\|^2}$$

- Show that  $\mathbf{K}$  is positive definite by computing its eigenvalues

```
n = 10;
X=randn(n,5);           % matrix of vectors
K = zeros(n,n);         % define holder for Gram
for i = 1:n
    x = X(i,:);
    for j = 1:n
        y = X(j,:);
        K(i,j) = exp(-norm(x-y)^2); % define elements of Gram matrix
    endfor
endfor

K           % Gram matrix
eig(K)      % Eigenvalues should all be non-negative
```

## 5 Answers

### 5.1 Quadratic Kernel

- This is just straightforward algebra

$$\begin{aligned}\phi^T(\mathbf{x})\phi(\mathbf{y}) &= x_1^2 y_1^2 + x_2^2 y_2^2 + x_3^2 y_3^2 + 2x_1 x_2 y_1 y_2 + 2x_1 x_3 y_1 y_3 + 2x_2 x_3 y_2 y_3 \\ &= (x_1 y_1 + x_2 y_2 + x_3 y_3)^2 = (\mathbf{x}^T \mathbf{y})^2\end{aligned}$$

- In the lecture notes we did the 2-d case
- Note that the more general polynomial kernel is

$$K_p(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x}^T \mathbf{y})^p$$

– this is more commonly used as it incorporates the lower dimensional polynomial kernels