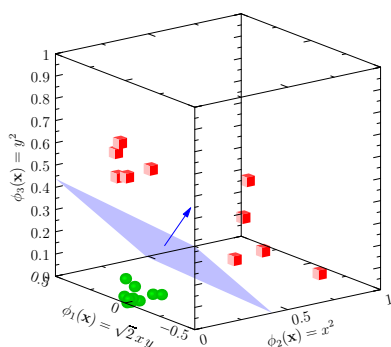


## Advanced Machine Learning

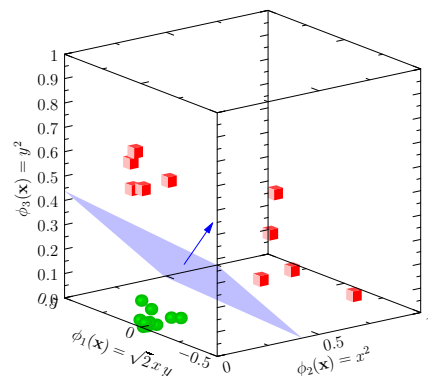
### Kernel Trick



The Kernel Trick, SVMs, Regression

## Outline

1. The Kernel Trick
2. Positive Semi-Definite Kernels
3. Kernel Properties
4. Beyond Classification



## SVM Kernels

- SVM Kernels are functions of two variables that can be factorised

$$K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle = \sum_i \phi_i(\mathbf{x}) \phi_i(\mathbf{y})$$

- where  $\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots)^T$  and  $\phi_i(\mathbf{x})$  are real valued functions of  $\mathbf{x}$
- $K(\mathbf{x}, \mathbf{y})$  will be positive semi-definite (because it is an inner-product)
- Furthermore, any positive semi-definite function will factorise
- This factorisation is not always obvious (we return to this later)

## Dual Form

- Recall that the dual problem for an SVM is

$$\max_{\alpha} \sum_{k=1}^m \alpha_k - \frac{1}{2} \sum_{k,l=1}^m \alpha_k \alpha_l y_k y_l \langle \phi(\mathbf{x}_k), \phi(\mathbf{x}_l) \rangle$$

- subject to  $\sum_{k=1}^m y_k \alpha_k = 0$  and  $0 \leq \alpha_k (\leq C)$
- But since  $K(\mathbf{x}_k, \mathbf{x}_l) = \langle \phi(\mathbf{x}_k), \phi(\mathbf{x}_l) \rangle$  the dual problem becomes

$$\max_{\alpha} \sum_{k=1}^m \alpha_k - \frac{1}{2} \sum_{k,l=1}^m \alpha_k \alpha_l y_k y_l K(\mathbf{x}_k, \mathbf{x}_l)$$

- This is the **kernel trick**—we never have to compute  $\phi(\mathbf{x})$

## Classifying New Data

- Having trained the SVM we now have to use it
- Given a new input  $\mathbf{x}$  we decide on the class

$$y = \text{sgn}(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle - b) \quad \text{but} \quad \mathbf{w} = \sum_{k=1}^m \alpha_k y_k \phi(\mathbf{x}_k)$$

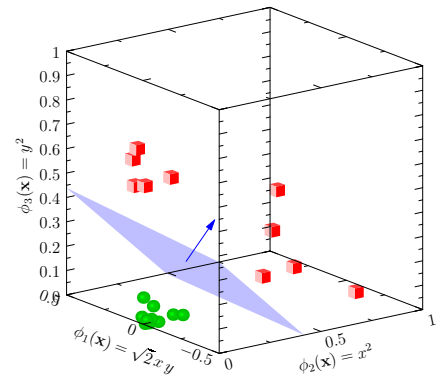
- In the dual representation this becomes

$$\text{sgn} \left( \sum_{k=1}^m \alpha_k y_k K(\mathbf{x}_k, \mathbf{x}) - b \right)$$

where we only need to sum over the non-zero  $\alpha_k$  (i.e. the support vectors SVs)

## Outline

1. The Kernel Trick
2. **Positive Semi-Definite Kernels**
3. Kernel Properties
4. Beyond Classification



## Recap on Eigen Systems

- Recall for a symmetric ( $n \times n$ ) matrix  $\mathbf{M}$  an eigenvector,  $\mathbf{v}$

$$\mathbf{M}\mathbf{v} = \lambda\mathbf{v}$$

- There are  $n$  independent eigenvectors  $\mathbf{v}^{(i)}$  with real eigenvalues  $\lambda^{(i)}$
- The eigenvectors are orthogonal so that  $\mathbf{v}^{(i)\top}\mathbf{v}^{(j)} = 0$  if  $i \neq j$
- Forming a matrix of eigenvectors  $\mathbf{V} = (\mathbf{v}^{(1)}, \mathbf{v}^{(2)}, \dots, \mathbf{v}^{(n)})$  the matrix satisfies

$$\mathbf{V}^\top \mathbf{V} = \mathbf{V} \mathbf{V}^\top = \mathbf{I}$$

- Such matrices are said to be orthogonal

## Eigen Decomposition

- From the eigenvalue equation  $\mathbf{M}\mathbf{v}^{(k)} = \lambda^{(k)}\mathbf{v}^{(k)}$

$$\mathbf{M}\mathbf{V} = \mathbf{V}\mathbf{\Lambda} \quad \text{where} \quad \mathbf{\Lambda} = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n \end{pmatrix}$$

- Multiplying on the right by  $\mathbf{V}^\top$  we get

$$\mathbf{M} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top = \sum_{k=1}^n \lambda^{(k)} \mathbf{v}^{(k)} \mathbf{v}^{(k)\top}$$

Or

$$M_{ij} = \sum_{k=1}^n \lambda^{(k)} v_i^{(k)} v_j^{(k)} = \sum_{k=1}^n u_i^{(k)} u_j^{(k)} = \langle \mathbf{u}_i, \mathbf{u}_j \rangle$$

$$u_i^{(k)} = \sqrt{\lambda^{(k)}} v_i^{(k)}$$

## Eigenfunctions

- By analogy for a symmetric function of two variables we can define an *eigenfunction*

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

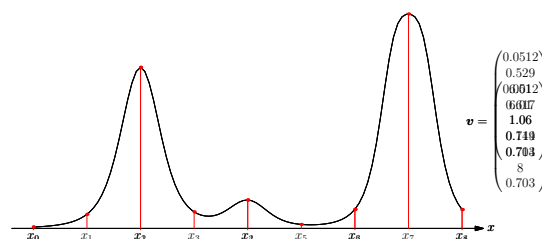
- In general there will be a denumerable set of eigenfunctions  $\psi^{(k)}(\mathbf{x})$  where

$$K(\mathbf{x}, \mathbf{y}) = \sum_k \lambda^{(k)} \psi^{(k)}(\mathbf{x}) \psi^{(k)}(\mathbf{y})$$

- This is known as Mercer's theorem

## Limit Process

- Consider sampling a function at a set of points

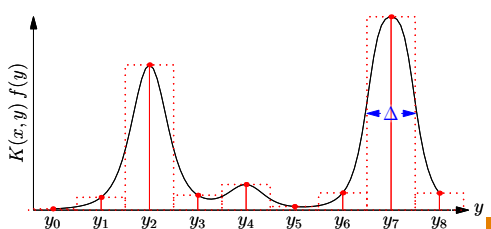


- In the limit where the number of sample points goes to infinity the vector more closely approximates a function
- Instead of the indices being numbers we use  $k \leftarrow x_k$

## Linear Operators

- Recall a linear function  $\mathcal{T}[f(x)]$  can be represented by a kernel

$$\mathcal{T}[f(x)] = \int_{y \in \mathcal{I}} K(x, y) f(y) dy \approx \Delta \sum_{j=1}^n K(x, y_j) f(y_j)$$



This is just a matrix equation with  $M_{ij} = \Delta K(x_i, y_j)$

## SVM Kernels

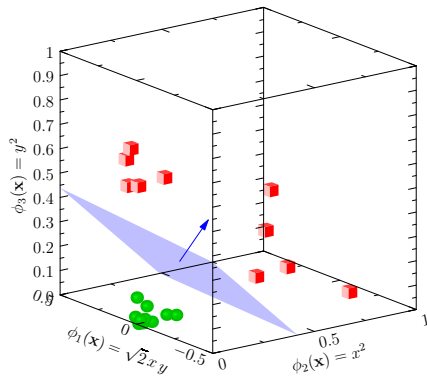
- If we define  $\phi^{(k)}(\mathbf{x}) = \sqrt{\lambda^{(k)}} \psi^{(k)}(\mathbf{x})$  then

$$K(\mathbf{x}, \mathbf{y}) = \sum_k \lambda^{(k)} \psi^{(k)}(\mathbf{x}) \psi^{(k)}(\mathbf{y}) = \sum_k \phi^{(k)}(\mathbf{x}) \phi^{(k)}(\mathbf{y})$$

- This is the definition of a SVM kernel we started with
- Note that for  $\phi^{(k)}(\mathbf{x})$  to be real  $\lambda^{(k)} \geq 0$  for all  $k$
- If  $\lambda^{(k)} < 0$  then  $\langle \phi(\mathbf{x}), \phi(\mathbf{x}) \rangle = \|\phi(\mathbf{x})\|^2$  might be negative and “distance” between points in the extended feature space can be negative
- If we use a kernel that isn't positive semi-definite then the Hessian of the dual objective function will not be negative semi-definite and there will be a maximum where  $\alpha$  diverges

## Outline

1. The Kernel Trick
2. Positive Semi-Definite Kernels
3. Kernel Properties
4. Beyond Classification



## Positive Semi-Definite Kernels

- Kernels (or matrices) that have eigenvalues  $\lambda^{(k)} \geq 0$  are called positive semi-definite
- (If the eigenvalues are strictly positive  $\lambda^{(k)} > 0$  the kernels or matrices are called positive definite)
- Positive semi-definite kernels can always be decomposed into a sum of real functions

$$K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$$

## Properties of Positive Semi-Definiteness

- Since

$$K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$$

- An immediate consequence is that for any function  $f(\mathbf{x})$

$$\begin{aligned} \int f(\mathbf{x}) K(\mathbf{x}, \mathbf{y}) f(\mathbf{y}) d\mathbf{x} d\mathbf{y} &= \int f(\mathbf{x}) \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle f(\mathbf{y}) d\mathbf{x} d\mathbf{y} \\ &= \left\langle \int f(\mathbf{x}) \phi(\mathbf{x}) d\mathbf{x}, \int f(\mathbf{y}) \phi(\mathbf{y}) d\mathbf{y} \right\rangle \\ &= \left\| \int f(\mathbf{x}) \phi(\mathbf{x}) d\mathbf{x} \right\|^2 \geq 0 \end{aligned}$$

## Positive Semi-Definiteness

- The following statements are equivalent
  - ★  $K(\mathbf{x}, \mathbf{y})$  is positive semi-definite (written  $K(\mathbf{x}, \mathbf{y}) \succeq 0$ )
  - ★ The eigenvalues of  $K(\mathbf{x}, \mathbf{y})$  are non-negative
  - ★ The kernel can be written

$$K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$$

where the  $\phi^{(k)}(\mathbf{x})$ 's are real functions

- ★ For any real function  $f(\mathbf{x})$

$$\int f(\mathbf{x}) K(\mathbf{x}, \mathbf{y}) f(\mathbf{y}) d\mathbf{x} d\mathbf{y} \geq 0$$

## Adding Kernels

- We can construct SVM kernels from other kernels
- If  $K_1(\mathbf{x}, \mathbf{y})$  and  $K_2(\mathbf{x}, \mathbf{y})$  are valid kernels then so is  $K_3(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x}, \mathbf{y}) + K_2(\mathbf{x}, \mathbf{y})$

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

- If  $K(\mathbf{x}, \mathbf{y})$  is a valid kernel so is  $cK(\mathbf{x}, \mathbf{y})$  for  $c > 0$

## Exponentiating Kernels

- If  $K(\mathbf{x}, \mathbf{y})$  is a valid kernel so is  $K^n(\mathbf{x}, \mathbf{y})$  (by induction)
  - ★ Assume  $K(\mathbf{x}, \mathbf{y}) \succeq 0$  this satisfies base case
  - ★ If  $K^{n-1}(\mathbf{x}, \mathbf{y}) \succeq 0$  then

$$K^n(\mathbf{x}, \mathbf{y}) = K^{n-1}(\mathbf{x}, \mathbf{y}) K(\mathbf{x}, \mathbf{y}) \succeq 0$$

- and  $\exp(K(\mathbf{x}, \mathbf{y}))$  is also a valid kernel since

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

but each term in the sum is a kernel

## Product of Kernels

- If  $K_1(\mathbf{x}, \mathbf{y})$  and  $K_2(\mathbf{x}, \mathbf{y})$  are valid kernels then so is  $K_3(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x}, \mathbf{y}) K_2(\mathbf{x}, \mathbf{y})$

- Writing

$$K_1(\mathbf{x}, \mathbf{y}) = \sum_i \phi_i^{(1)}(\mathbf{x}) \phi_i^{(1)}(\mathbf{y}), \quad K_2(\mathbf{x}, \mathbf{y}) = \sum_j \phi_j^{(2)}(\mathbf{x}) \phi_j^{(2)}(\mathbf{y})$$

then

$$\begin{aligned} K_3(\mathbf{x}, \mathbf{y}) &= \sum_{i,j} \phi_i^{(1)}(\mathbf{x}) \phi_i^{(1)}(\mathbf{y}) \phi_j^{(2)}(\mathbf{x}) \phi_j^{(2)}(\mathbf{y}) \\ &= \sum_{i,j} \phi_{ij}^{(3)}(\mathbf{x}) \phi_{ij}^{(3)}(\mathbf{y}) = \langle \phi^{(3)}(\mathbf{x}), \phi^{(3)}(\mathbf{y}) \rangle \end{aligned}$$

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

## RBF Kernel

- Now  $\mathbf{x}^\top \mathbf{y} = \langle \mathbf{x}, \mathbf{y} \rangle$  is a valid kernel because it is an inner product of functions  $\phi(\mathbf{x}) = \mathbf{x}$
- For  $\gamma > 0$  we have  $2\gamma \mathbf{x}^\top \mathbf{y} \succeq 0$
- Thus  $\exp(2\gamma \mathbf{x}^\top \mathbf{y}) \succeq 0$
- If  $K(\mathbf{x}, \mathbf{y}) \succeq 0$  then  $g(\mathbf{x}) K(\mathbf{x}, \mathbf{y}) g(\mathbf{y}) \succeq 0$

$$\int f(\mathbf{x}) g(\mathbf{x}) K(\mathbf{x}, \mathbf{y}) g(\mathbf{y}) f(\mathbf{y}) d\mathbf{x} d\mathbf{y} = \int h(\mathbf{x}) K(\mathbf{x}, \mathbf{y}) h(\mathbf{y}) d\mathbf{x} d\mathbf{y} \succeq 0$$

where  $f(\mathbf{x}) g(\mathbf{x}) = h(\mathbf{x})$

$$e^{-\gamma \mathbf{x}^\top \mathbf{x}} e^{2\gamma \mathbf{x}^\top \mathbf{y}} e^{-\gamma \mathbf{y}^\top \mathbf{y}} = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|^2} \succeq 0$$

## Other Kernels

- The success of SVMs has meant that researchers try to increase the area of application
- The condition that a SVM kernel must be positive semi-definite is quite restrictive
- There has been a cottage industry of researchers finding smart kernels for solving complicated problems
- The key to finding new kernels is to use the properties of kernels to build more complicated kernels from simpler ones

## String Kernels

- One area where SVMs were very important is in document classification
- This requires comparing strings
- There are a large number of kernels developed to do this

## Spectrum Kernel

- A simple way to compare documents is to collect a histogram of all occurrences of substrings of length  $p$
- This is known as a  $p$ -spectrum
- A  $p$ -spectrum kernel counts the number of common substrings

$s = \text{statistics}$      $S_3(s) = \{\text{sta}, \text{tat}, \text{ati}, \text{tis}, \text{ist}, \text{sti}, \text{tic}, \text{ics}\}$   
 $t = \text{computation}$      $S_3(t) = \{\text{com}, \text{omp}, \text{mpu}, \text{put}, \text{uta}, \text{tat}, \text{ati}, \text{tio}, \text{ion}\}$

- $K(s, t) = 2$  ("tat" and "ati")

## All Subsequences Kernel

- A more sophisticated kernel is to count all of the common subsequences that occur in two documents
- Naively this would take an exponential amount of time to compute
- Using clever dynamic-programming techniques this can be done relatively efficiently
- This can even be extended to include sub-sequence matches with possible gaps between words

## Other Kernel Applications

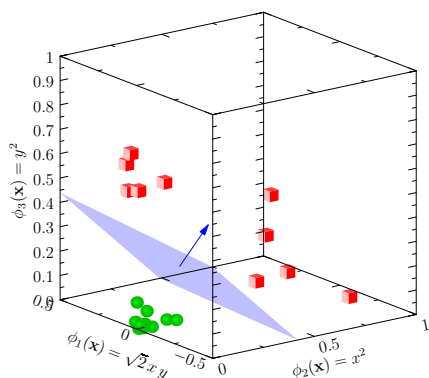
- String kernels for comparing subsequences are used in bioinformatics
- Kernels have been developed for comparing trees (e.g. for computer program evaluation, XML, etc.)
- Kernels have also been developed for comparing graphs (e.g. for comparing chemicals based on their molecular graph)

## Fisher Kernels

- In an attempt to build kernels that capture more domain knowledge, kernels are constructed from other learning machines
- An example of this are “Fisher kernels” whose features come from an Hidden Markov Model (HMM) trained on the data
- These tend to have better discriminative power than the underlying model (HMM), and has a better feature set than a SVM using a generic kernel

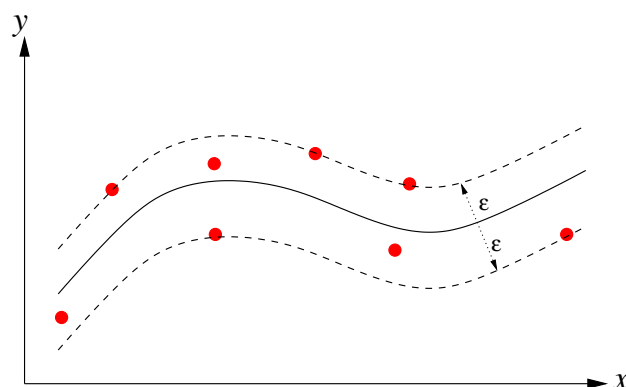
## Outline

1. The Kernel Trick
2. Positive Semi-Definite Kernels
3. Kernel Properties
4. Beyond Classification



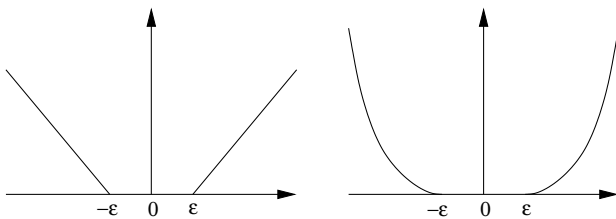
## Regression with Margins

- SVMs can be modified to perform regression



## Error Functions

- Can introduce slack variables with different errors



- This can be transformed to a quadratic programming problem

## Kernel Methods

- Kernel methods where we project into an extended feature space are used with other linear algorithms
  - ★ Kernel Fisher discriminant analysis (KFDA)
  - ★ Kernel principle component analysis (KPCA)
  - ★ Kernel canonical correlation analysis (KCCA)
  - ★ Gaussian Processes
- These are also extremely powerful machine learning algorithms

## Ridge Regression Using Kernels

- We can also solve regression problems without using margins
- To solve a regression problem once again the problem is set up as a quadratic programming problem

$$\min_{\mathbf{w}} \lambda \|\mathbf{w}\|^2 + \sum_{i=1}^m (y_i - \mathbf{w}^T \phi(x_i))^2$$

- the  $\|\mathbf{w}\|^2$  is a regularisation term
- By assuming  $\mathbf{w} = \sum_i \alpha_i \phi(x_i)$  we obtain a quadratic equation for the  $\alpha_i$ 's which we can solve

## Summary

- SVMs require a positive definite kernel function
- These can be built from simpler functions
- There was a cottage industry of people creating new kernels for different applications
- SVMs are just one example of a host of machine learning algorithms that
  - ★ use the kernel trick
  - ★ often use linear constraints
  - ★ tend to be convex optimisation problems