

Kernel Properties, positive semi-definiteness, string kernels

SVMs

- A linear SVM finds the maximal margin hyperplane for separating linear separable data
- We can increase the chances of the data being linearly separable by projecting the data into an **extended feature space**

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix} \rightarrow \phi(\mathbf{x}) = \begin{pmatrix} \phi_1(\mathbf{x}) \\ \phi_2(\mathbf{x}) \\ \vdots \\ \phi_r(\mathbf{x}) \end{pmatrix}$$

Kernel Trick

- If we define the kernel function as

$$K(\mathbf{x}, \mathbf{y}) = \phi^T(\mathbf{x}) \phi(\mathbf{y})$$

then

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

$$\hat{y} = \text{sgn} \left(\sum_{k \in SV} \alpha_k y_k K(\mathbf{x}_k, \mathbf{x}) - b \right)$$

- We only need to compute the kernel rather than $\phi(\mathbf{x})$

Eigen-Functions

- In analogy to eigen-vectors we can define eigen-functions of a function of two variables

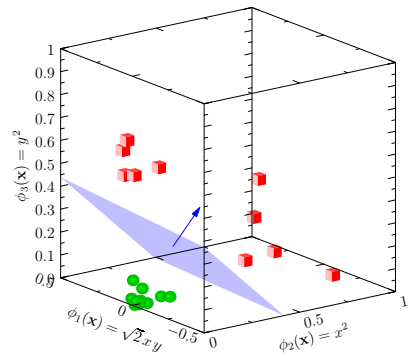
$$\int K(\mathbf{x}, \mathbf{y}) \psi_k(\mathbf{y}) d\mathbf{y} = \lambda_k \psi_k(\mathbf{x})$$

$$\sum_j M_{ij} v_j^{(k)} = \lambda_k v_i^{(k)}$$

- The spatial coordinate \mathbf{x} plays the same role as the index i
- We can decompose a kernel into a sum of its eigen-functions

$$K(\mathbf{x}, \mathbf{y}) = \sum_i \lambda_i \psi_i(\mathbf{x}) \psi_i(\mathbf{y}) \quad \text{c.f.} \quad \mathbf{M} = \sum_{i=1} \lambda_i \mathbf{v}_i \mathbf{v}_i^T$$

- Recap
- Positive Semi-Definite Kernels
- Training SVMs
- Beyond Classification



Dual Form

- Finding the maximum margin hyperplane is equivalent to solving the quadratic programming problem

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

subject to $\alpha_k \geq 0$ and $\sum_k y_k \alpha_k = 0$

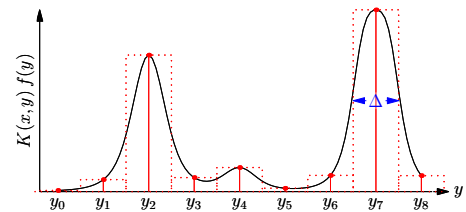
- This uses the data set $\{(\mathbf{x}_k, y_k) | k = 1, \dots, m\}$ to learn a set of α_k 's
- To classify new data we get a class prediction

$$\hat{y} = \text{sgn} \left(\sum_{k \in SV} \alpha_k y_k \phi^T(\mathbf{x}_k) \phi(\mathbf{x}) - b \right)$$

Kernels and Matrices

- A linear transformation $\mathcal{T}[f(x)]$ 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 $M_{ij} = \Delta K(x_i, y_j)$

Mercer's Theorem

- Mercer tells us that for any symmetric kernel function

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

- If $\lambda_i \geq 0$ for all i then we can define $\phi_i(\mathbf{x}) = \sqrt{\lambda_i} \psi_i(\mathbf{x})$
- And

$$K(\mathbf{x}, \mathbf{y}) = \sum_i \phi_i(\mathbf{x}) \phi_i(\mathbf{y}) = \phi^T(\mathbf{x}) \phi(\mathbf{y})$$

- That is, any positive semi-definite symmetric function of two variables is a valid kernel function!

- If we used any old kernel function with some negative eigenvalues then there can be a projection $x \rightarrow \phi(x)$ such that

$$\phi(x)^T \phi(x) < 0$$

(e.g. if $\phi(x)$ was an eigenvector with negative eigenvalue)

- We are no longer in a space with Euclidean geometry
- Maximum margins are meaningless
- Must use positive semi-definite kernels

Positive Semi-Definite Kernels

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

$$K(x, y) = \sum_i \phi_i(x) \phi_i(y)$$

Positive Semi-Definiteness

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

$$K(x, y) = \sum_i \phi_i(x) \phi_i(y)$$

where $\phi(x)$ are real functions

- For any real function $f(x)$

$$\int f(x) K(x, y) f(y) dx dy \geq 0$$

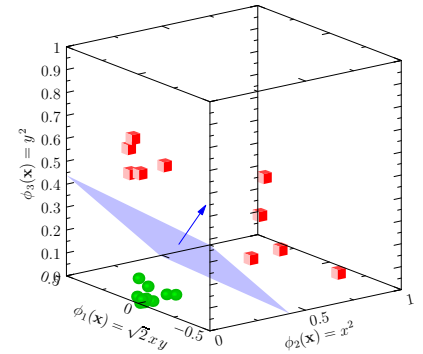
Product of Kernels

- If $K_1(x, y)$ and $K_2(x, y)$ are valid kernels then so is $K_3(x, y) = K_1(x, y) K_2(x, y)$
- Writing $K_1(x, y) = \sum_i \phi_i^1(x) \phi_i^1(y)$ and $K_2(x, y) = \sum_j \phi_j^2(x) \phi_j^2(y)$ then

$$\begin{aligned} K_3(x, y) &= K_1(x, y) K_2(x, y) = \sum_{i,j} \phi_i^1(x) \phi_i^1(y) \phi_j^2(x) \phi_j^2(y) \\ &= \sum_{i,j} (\phi_i^1(x) \phi_j^2(x)) (\phi_i^1(y) \phi_j^2(y)) \\ &= \sum_{i,j} \phi_{ij}^3(x) \phi_{ij}^3(y) \end{aligned}$$

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

- Recap
- Positive Semi-Definite Kernels**
- Training SVMs
- Beyond Classification



Properties of Positive Semi-Definiteness

- Since

$$K(x, y) = \sum_i \phi_i(x) \phi_i(y)$$

- An immediate consequence is that for any function $f(x)$

$$\begin{aligned} \int f(x) K(x, y) f(y) dx dy &= \int f(x) \sum_i \phi_i(x) \phi_i(y) f(y) dx dy \\ &= \sum_i \left(\int f(x) \phi_i(x) dx \right)^2 \geq 0 \end{aligned}$$

Adding Kernels

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

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

- If $K(x, y)$ is a valid kernel so is $c K(x, y)$ for $c > 0$

Exponentiating Kernels

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

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

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

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

but each term in the sum is a kernel

- Now $\mathbf{x}^T \mathbf{y}$ is a valid kernel because it is of the form $\sum_i \phi_i(\mathbf{x}) \phi_i(\mathbf{y})$ where $\phi_i(\mathbf{x}) = x_i$
- For $\gamma > 0$ we have $2\gamma \mathbf{x}^T \mathbf{y} \succeq 0$
- Thus $\exp(2\gamma \mathbf{x}^T \mathbf{y}) \succeq 0$
- Since $\exp(-\gamma \mathbf{x}^T \mathbf{x})$ and $\exp(-\gamma \mathbf{y}^T \mathbf{y})$ are positive numbers

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

- This is the RBF or Gaussian kernel

String Kernels

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

All Subsequences Kernel

- A more sophisticated kernel is to count all of the common subsequences that occur in two documents
- Naively this would take a huge 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

Fisher Kernels

- In an attempt to build kernels that capture more domain knowledge, kernels are constructed from other learning machines
- “Fisher kernels” can be constructed from features coming from generative models (e.g. a Hidden Markov Model (HMM) trained on biological 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

- 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 an industry of research 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

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

$$\begin{aligned} s = \text{statistics} \quad \mathcal{S}_3(s) &= \{\text{sta, tat, ati, tis, ist, sti, tic, ics}\} \\ t = \text{computation} \quad \mathcal{S}_3(t) &= \{\text{com, omp, mpu, put, uta, tat, ati, tio, ion}\} \end{aligned}$$

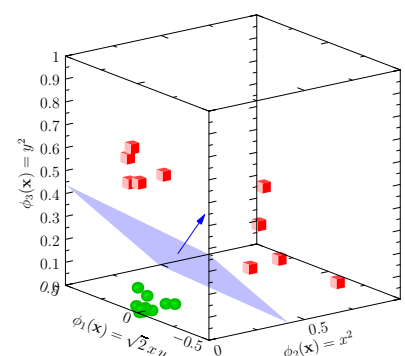
- $K(s, t) = 2$ (“tat” and “ati”)

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)

Outline

1. Recap
2. Positive Semi-Definite Kernels
3. Training SVMs
4. Beyond Classification



- The dual problem is

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

subject to $\alpha_k \geq 0$ and $\sum_k y_k \alpha_k = 0$

- If we allow slack variables with a constraint $C \sum_k \xi_k$ then get the same problem with

$$0 \leq \alpha_k \leq C \quad \forall k = 1, 2, \dots, P$$

Time Complexity of SVMs

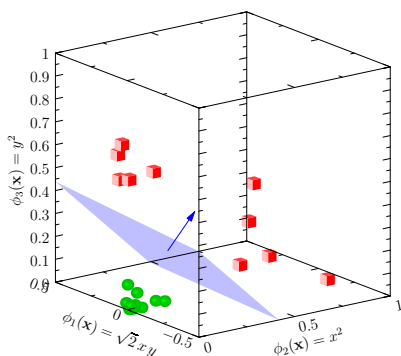
- SVMs have good generalisation performance often beating MLP or RBFs
- They have a unique classification boundary unlike MLPs which can find local optima
- The time complexity for training is $O(m^3)$ where m is the number of training patterns
- It can be too slow if $m \gg 1000$

Sequential Minimal Optimisation

- One of the most efficient techniques for training SVMs is *Sequential Minimal Optimisation* or SMO
- This takes two Lagrange multipliers α_i and α_j and adjusts them to maximise the dual objective function
- This is very quick as it can be done in closed form
- Note that because $\sum_k y_k \alpha_k = 0$ we have to change at least two variables at the same time
- A heuristic is used to choose the best pair of α 's to optimise
- Run until close to the optimum

Outline

- Recap
- Positive Semi-Definite Kernels
- Training SVMs
- Beyond Classification



- Traditional quadratic programming solvers start from a feasible solutions (usually found using linear programming)
- Takes the current set of constraints that are exactly satisfied as the active set
- Optimises with respect to the active set taken as equality constraints
- Moves towards the new optimum as far as possible (so that none of the non-active constraints is broken)
- If any of the Lagrange multipliers are negative, it drops the constraints

Chunking

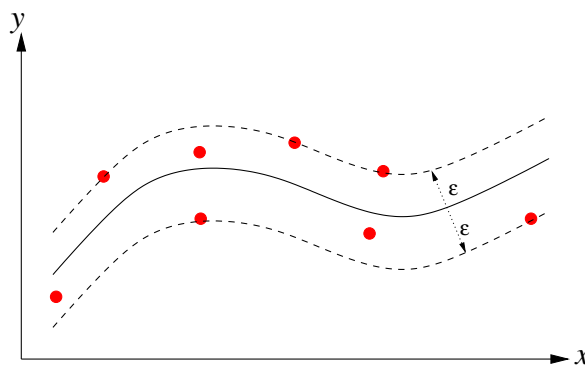
- Chunking is an attempt to find an approximation for the maximum margin hyperplane by working on "chunks" of the dataset
- The algorithm considers an *chunk* of data at a time
- Initially we run an SVM on the first chunk of data
- The support vectors for the chunk are retained while the rest of the data in the chunk is discarded
- The support vectors together with the next set of data is considered

Multi-class Classification

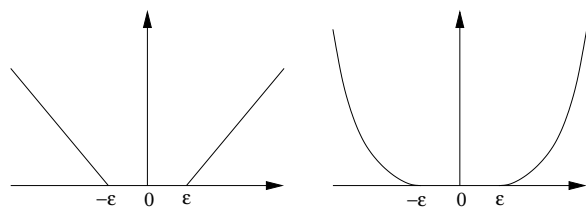
- SVMs are by nature a binary classifier
- There are a number of strategies to make them multi-class, two frequent strategies
 - ★ Train $|\mathcal{C}|$ one-versus-all classifiers and choose best
 - ★ Train $|\mathcal{C}|(|\mathcal{C}| - 1)/2$ one-versus-one classifiers and vote for best
- More elegant, but slightly more complicated alternatives exist involving using the class label as a feature

Regression with Margins

- SVMs can be modified to perform regression



- 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 is also used with algorithms
 - ★ Fisher discriminant analysis
 - ★ Principle component analysis
 - ★ Canonical correlation analysis
 - ★ Gaussian Processes
- These are also extremely power machine learning algorithms

- 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

$$\min_{\mathbf{w}} \lambda \|\mathbf{w}\|^2 + \sum_{k=1}^m (y_k - \mathbf{w}^T \phi(\mathbf{x}_k))^2$$

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

Summary

- SVMs require a positive definite kernel function
- These can be built from simpler function
- There is an important industry of people creating new kernels for different application
- SVMs can be slow for very large datasets, but there are approximation methods to get around this
- There are lots of good SVM libraries, but care is need using them (normalising inputs and tuning parameters)
- Even when you understand all the mathematics. . . they are still magic