

Introduction to Kernel methods

Chiranjib Bhattacharyya

Machine Learning lab

Dept of CSA, IISc

`chiru@csa.iisc.ernet.in`

`http://drona.csa.iisc.ernet.in/~chiru`

13th June, 2011

Agenda

- 1 Kernel Trick
 - SVM and Non-linear Classification
- 2 Definition of Kernel functions
- 3 Kernels and Hilbert Spaces
 - RKHS, Representer theorem etc
- 4 Kernels on Different kinds of data
 - Kernels on Strings
 - Kernels on generative models
 - Kernels on Graphs
- 5 Advanced Topics
 - Multiple Kernel Learning
 - Learning Kernels from Similarity functions

Introduction

- **Data:** Protein structures, Images, videos,
Algorithms: Classifiers, PCA

Most algorithms may require vectorial description which maynot be easily available. Kernels can help bridge the gap.

- Kernels are similarity measures
- Kernels can help integrate different sources of data

PART 1: KERNEL TRICK

The problem of classification

Let $D = \{(x_i, y_i) | i = 1, \dots, m\}$ be the training dataset. where y_i is the class label of an observation x_i

Assume that $y \in \{-1, 1\}$.

A classifier $f : \mathcal{X} \rightarrow \{-1, 1\}$. A linear classifier is of the form

$$f(x) = \text{sign}(w^\top x + b)$$

Learning classifiers on a given dataset D which has good *generalization* abilities is one of the important problems in Machine Learning.

Review of C-SVM

$$\min_{\mathbf{w}, b} C \sum_{i=1}^m \max(1 - y_i(\mathbf{w}^\top x_i + b), 0) + \frac{1}{2} \|\mathbf{w}\|^2$$

C-SVM formulation

$$\begin{aligned} \text{maximize}_{\alpha} & -\frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j x_i^\top x_j + \sum_{i=1}^m \alpha_i \\ \text{subject to} & 0 \leq \alpha_i, \sum_i \alpha_i y_i = 0 \end{aligned}$$

At optimality $\mathbf{w} = \sum_{i=1}^m \alpha_i y_i x_i$

$$f(x) = \text{sign}\left(\sum_{i=1}^m \alpha_i y_i x_i^\top x + b\right)$$

C-SVM in feature spaces

Let us work with a feature map, $\Phi(x)$ Solving a SVM problem in feature space turns out to be

$$\text{maximize}_{\alpha} -\frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j \Phi(x_i)^\top \Phi(x_j) + \sum_{i=1}^m \alpha_i$$

$$\text{subject to } 0 \leq \alpha_i, \sum_i \alpha_i y_i = 0$$

and our classifier is

$$f(x) = \text{sign}\left(\sum_{i=1}^m \alpha_i y_i \Phi(x_i)^\top \Phi(x) + b\right)$$

The **dot product** between any pair of examples computed in the feature space be denoted by

$$K(x, z) = \Phi(x)^\top \Phi(z)$$

C-SVM in feature spaces

Let us work with a feature map, $\Phi(x)$ Solving a SVM problem in feature space turns out to be

$$\text{maximize}_{\alpha} -\frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j K(x_i, x_j) + \sum_{i=1}^m \alpha_i$$

$$\text{subject to } 0 \leq \alpha_i, \sum_i \alpha_i y_i = 0$$

and our classifier is

$$f(x) = \text{sign}\left(\sum_{i=1}^m \alpha_i y_i K(x_i, x) + b\right)$$

The **dot product** between any pair of examples computed in the feature space be denoted by

$$K(x, z) = \Phi(x)^\top \Phi(z)$$

Kernel Fisher Discriminant(Mika et al. 1999)

Let $X = [x_1, \dots, x_{m_+}]$ and $Z = [z_1, \dots, z_{m_-}]$ where $x_i, z_j \in \mathbb{R}^d$

$$\max_w \frac{\{w^\top(\mu_1 - \mu_2)\}^2}{w^\top S_w w}$$

$$\mu_1 = \frac{1}{m_+} \sum_i x_i = \frac{1}{m_+} X e \quad \mu_2 = \frac{1}{m_-} \sum_i z_i = \frac{1}{m_-} Z e$$

(e is a vector of ones)

$$S_w = \frac{1}{m_+} \sum_{i=1}^{m_+} (x_i - \mu_+)(x_i - \mu_+)^\top + \frac{1}{m_-} \sum_{j=1}^{m_-} (z_j - \mu_-)(z_j - \mu_-)^\top$$

$\forall w \in \mathbb{R}^d$ there exists unique w_\perp such that

$$w = [X \ Z] \alpha + w_\perp \quad w_\perp^\top x_i = w_\perp^\top z_j = 0$$

Kernel Fisher Discriminant(Mika et al. 1999)

$$w^\top \mu_+ = \alpha^\top [X^\top X Z^\top X]e \quad w^\top \mu_- = \alpha^\top [X^\top Z Z^\top Z]e$$

$$(X^\top Z)_{ij} = x_i^\top z_j$$

can be **kernelized**, i.e. $x_i^\top z_j$ can be replaced by kernel functions $K(x_i, z_j)$.

Kernel Fisher Discriminant(Mika et al. 1999)

$$w^\top \mu_+ = \alpha^\top [X^\top X Z^\top X] e \quad w^\top \mu_- = \alpha^\top [X^\top Z Z^\top Z] e$$

$$(X^\top Z)_{ij} = x_i^\top z_j$$

can be **kernelized**, i.e. $x_i^\top z_j$ can be replaced by kernel functions $K(x_i, z_j)$. It is possible to compute Fisher Discriminant in feature space

$$\max_{\alpha} \frac{\{\alpha^\top (k_1 - k_2)\}^2}{\alpha^\top C \alpha}$$

where k_1, k_2 and C are suitably defined.

Principal Component Analysis(PCA)

Let $X = [x_1, \dots, x_m]$. Finding directions of maximum variance can be very informative(see Jolliffe 2002).

Using sample covariance the direction of maximum variance, v , is given by

$$\frac{1}{m}XX^T v = \lambda v$$

(assuming that $Xe = 0$)

Let $\alpha = X^T v$ then $v = \frac{1}{\lambda} \sum_{i=1}^m \alpha_i x_i$ indicating that v is in the linear span of $\{x_1, \dots, x_m\}$.

$$\frac{1}{m}XX^T X\alpha = \lambda X\alpha$$

leading to the following eigenvalue problem

$$\frac{1}{m}\mathbf{K}\alpha = \lambda \alpha$$

where $(\mathbf{K})_{ij} = (X^T X)_{ij} = x_i^T x_j$. Replacing $x_i^T x_j$ by $\Phi(x_i)^T \Phi(x_j)$ one can do PCA in feature spaces (Scholkopf et al. 1996)

An example

- Let $x \in \mathbb{R}^2$ and $\Phi(x) = [x_1^2 \ x_2^2 \ \sqrt{2}x_1x_2]^\top$

$$K(x, z) = \Phi(x)^\top \Phi(z) = x_1^2 z_1^2 + 2x_1 x_2 z_1 z_2 + x_2^2 z_2^2 = \langle x, z \rangle^2$$

- If $K(x, z) = (x^\top z)^r$ is a dot product in a $\binom{d+r-1}{r}$ feature space corresponding to $x, z \in \mathbb{R}^d$.
- If $d = 256, r = 4$, the feature space size is 6,35,376.
- However if we know K one can still solve the SVM formulation without explicitly evaluating Φ

In the sequel

- conditions on K will be discussed
- K for various data-types
- Advanced topics in Kernel design

Kernel function

Kernel function

$K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a Kernel function if

$$K(x, z) = K(z, x) \text{ symmetric}$$

K is positive semidefinite, i.e. $\forall n, x_1, \dots, x_n \in \mathcal{X}$,
the matrix $K_{ij} = K(x_i, x_j)$ is psd

Recall that a $\mathbf{K} \in \mathbb{R}^{d \times d}$ is psd if $u^\top \mathbf{K} u \geq 0$ for all $u \in \mathbb{R}^d$.

Examples of Kernel function

$$K(x, z) = \phi(x)^\top \phi(z) \text{ where } \phi : \mathcal{E} \rightarrow \mathbb{R}^d$$

K is symmetric i.e. $K(x, z) = K(z, x)$

Examples of Kernel function

$$K(x, z) = \phi(x)^\top \phi(z) \text{ where } \phi : \mathcal{E} \rightarrow \mathbb{R}^d$$

K is symmetric i.e. $K(x, z) = K(z, x)$

Positive Semidefinite:

Let $D = \{x_1, x_2, \dots, x_n\}$ be set of arbitrarily chosen n elements of \mathcal{E} .

Define

$$\mathbf{K}_{ij} = \phi(x_i)^\top \phi(x_j)$$

For any $u \in \mathbb{R}^n$ it is straightforward to see that

$$u^\top \mathbf{K} u = \|\phi(D)u\|_2^2 \geq 0 \quad \phi(D) = [\phi(x_1), \dots, \phi(x_n)]$$

Examples of Kernel functions

$$K(x, z) = x^\top z$$

$$K(x, z) = (x^\top z)^r$$

$$K(x, z) = e^{-\gamma \|x-z\|^2}$$

$$\Phi(x) = x$$

$$\Phi_{t_1 t_2 \dots t_d}(x) = \sqrt{\frac{r!}{t_1! t_2! \dots t_d!}} x_1^{t_1} x_2^{t_2} \dots x_d^{t_d}$$
$$\sum_{i=1}^d t_i = r$$

Kernel Construction

Let K_1 and K_2 be two valid kernels.

$$K(x, y) = \phi(x)^\top \phi(y)$$

$$K(u, v) = K_1(u, v)K_2(u, v)$$

$$K = \alpha K_1 + \beta K_2 \quad \alpha, \beta \geq 0$$

$$\hat{K}(x, y) = \frac{K(x, y)}{\sqrt{K(x, x)}\sqrt{K(y, y)}}$$

Kernel Construction

Let K_1 and K_2 be two valid kernels.

$$K(x, y) = \phi(x)^\top \phi(y)$$

$$K(u, v) = K_1(u, v)K_2(u, v)$$

$$K = \alpha K_1 + \beta K_2 \quad \alpha, \beta \geq 0$$

$$\hat{K}(x, y) = \frac{K(x, y)}{\sqrt{K(x, x)}\sqrt{K(y, y)}}$$

$$K(x, y) = x^\top y$$

$$K(x, y) = (x^\top y)^i$$

$$K(x, y) = \lim_{N \rightarrow \infty} \sum_{i=0}^N \frac{(x^\top y)^i}{i!} = e^{x^\top y}$$

$$\hat{K}(x, y) = e^{-\frac{1}{2}\|x-y\|^2}$$

Kernel function and feature map

A theorem due to **Mercer** guarantees a feature map for symmetric, psd kernel functions.

Loosely stated

For a symmetric kernel $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, there exists an expansion $K(x, z) = \Phi(x)^\top \Phi(z)$ iff

$$\int_{\mathcal{X}} g(x)g(z)K(x, z)dx dz \geq 0$$

PART 3: Kernels and Hilbert spaces

What is a Dot product(aka Inner Product)

Let \mathcal{X} be a vector space.

What is a Dot product

$$\text{Symmetry } \langle u, v \rangle = \langle v, u \rangle \quad u, v \in \mathcal{X}$$

$$\text{Bilinear } \langle \alpha u + \beta v, w \rangle = \alpha \langle u, w \rangle + \beta \langle v, w \rangle \quad u, v, w, \in \mathcal{X}$$

$$\text{Positive Semidefinite } \langle u, u \rangle \geq 0 \quad u \in \mathcal{X}$$

$$\langle u, u \rangle = 0 \text{ iff } u = 0$$

Norm

$$\|x\| = \sqrt{\langle x, x \rangle}$$

$$\|x\| = 0 \implies x = 0$$

Examples of Dot products

$$\mathcal{X} = \mathbb{R}^n, \langle u, v \rangle = u^\top v$$

$$\mathcal{X} = \mathbb{R}^n, \langle u, v \rangle = \sum_{i=1}^n \lambda_i u_i v_i \quad \lambda_i \geq 0$$

$$\mathcal{X} = L_2(X) = \left\{ f : \int_X f(x)^2 dx < \infty \right\}$$

$$f, g \in \mathcal{X} \quad \langle f, g \rangle = \int_X f(x)g(x)dx$$

Cauchy Schwartz inequality

Cauchy Schwartz inequality

Let \mathcal{X} be an *inner product space*.

$$|\langle x, y \rangle| \leq \|x\| \|y\| \quad \forall x, y \in \mathcal{X}$$

and equality holds iff $x = \alpha z$ for some scalar α

Proof: $\forall \alpha \in \mathbb{R} \quad \|x - \alpha z\|^2 \geq 0$

$$\|x\|^2 - 2\alpha \langle x, z \rangle + \alpha^2 \|z\|^2 \geq 0 \quad \forall \alpha$$

Let $\alpha = \frac{\langle x, z \rangle}{\|z\|^2}$ and the inequality follows by taking square roots. The claim about equality follows from the definition of norm.

Hilbert Space: Basic facts

Defn: A Inner product space $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ is a **Hilbert Space** if it is separable and complete.

We will denote the norm as $\|\cdot\|_{\mathcal{H}}$. The **orthogonal complement** of M , where $M \subset \mathcal{H}$ be a subspace of \mathcal{H} is defined as

$$M^{\perp} = \{z | \langle x, z \rangle_{\mathcal{H}} = 0, \forall x \in M\}$$

Hilbert space Projection theorem

Let M be a subspace of Hilbert space $\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}}$. For every $x \in \mathcal{H}$ the following holds

- There exists an unique $\Pi_M(x) \in M$ such that
$$\Pi_M(x) = \operatorname{argmin}_{z \in M} \|x - z\|_{\mathcal{H}}$$
- $x - \Pi_M(x) \in M^{\perp}$ $\langle z, x - \Pi_M(x) \rangle_{\mathcal{H}} = 0 \forall z \in M$
- $\|x\|_{\mathcal{H}}^2 = \|\Pi_M(x)\|_{\mathcal{H}}^2 + \|y\|_{\mathcal{H}}^2$ where

$$x = \Pi_M(x) + y \text{ where } y \in M^{\perp}$$

Reproducing kernel Hilbert Space(RKHS)

Let K be any kernel function. Consider the following set

$$\mathcal{H} = \{f|f(\cdot) = \sum_{i=1}^m \alpha_i K(\cdot, x_i) \forall x_i \in \mathcal{X}, m \in \mathbb{N}\}$$

Dot product

For any $f, g \in \mathcal{H}$,

$$f(\cdot) = \sum_{i=1}^{m_1} \alpha_i K(\cdot, x_i), \quad g(\cdot) = \sum_{j=1}^{m_2} \beta_j K(\cdot, x_j)$$

$$\langle f, g \rangle_{\mathcal{H}} = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \alpha_i \beta_j K(x_i, x_j)$$

Is it a dot product?

Reproducing kernel Hilbert Space(RKHS)

As K is symmetric, $\langle f, g \rangle_{\mathcal{H}} = \langle g, f \rangle_{\mathcal{H}}$

$$\langle f(\cdot), f(\cdot) \rangle = \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j K(x_i, x_j)$$

Recall that \mathbf{K} is a psd matrix if K is kernel function and so

$$\langle f(\cdot), f(\cdot) \rangle_{\mathcal{H}} \geq 0$$

Reproducible Property

for any $f \in \mathcal{H}$

$$f(x) = \sum_{i=1}^m \alpha_i K(x, x_i) = \left\langle \sum_{i=1}^m \alpha_i K(\cdot, x_i), K(\cdot, x) \right\rangle = \langle f(\cdot), K(\cdot, x) \rangle$$

Applying C-S inequality

$$|f(x)| \leq \sqrt{\langle f, f \rangle_{\mathcal{H}}} \sqrt{K(x, x)}$$

Representer theorem

Representer theorem

Let K be a valid kernel defined on \mathcal{X} and \mathcal{H} be the corresponding RKHS. Let Ω be an increasing function. The optimization problem

$$\min_{g \in \mathcal{H}} G(g) = \sum_{i=1}^m l(g(x_i), y_i) + \Omega(\|g\|_{\mathcal{H}}^2)$$

is solved when $g^* = \sum_{i=1}^m \alpha_i K(\cdot, x_i)$

Representer theorem

Representer theorem

Let K be a valid kernel defined on \mathcal{X} and \mathcal{H} be the corresponding RKHS. Let Ω be an increasing function. The optimization problem

$$\min_{g \in \mathcal{H}} G(g) = \sum_{i=1}^m l(g(x_i), y_i) + \Omega(\|g\|_{\mathcal{H}}^2)$$

is solved when $g^* = \sum_{i=1}^m \alpha_i K(\cdot, x_i)$

Proof: Let $M = \{\sum_{i=1}^m \alpha_i K(\cdot, x_i) \mid i = 1, \dots, m\}$. Clearly M is a subspace of \mathcal{H} . Take any $g \in \mathcal{H}$.

$$\begin{aligned} g(x_i) &= \langle g, K(\cdot, x_i) \rangle = \langle g_M + g_{per}, K(\cdot, x_i) \rangle \\ &= \langle g_M, K(\cdot, x_i) \rangle + \langle g_{per}, K(\cdot, x_i) \rangle = \langle g_M, K(\cdot, x_i) \rangle = g_M(x_i) \end{aligned}$$

As Ω is an increasing function, $\Omega(\|g\|_{\mathcal{H}}^2) \geq \Omega(\|g_M\|_{\mathcal{H}}^2)$

Back to C-SVM formulation

Given a Kernel function K defined on \mathcal{X} one can create a RKHS

$$\mathcal{H} = \left\{ \sum_{i=1}^n \beta_i z_i \mid z_i \in \mathcal{X}, n \in \mathbb{N} \right\}$$

The C-SVM problem can be posed as

$$\min_{g \in \mathcal{H}, b \in \mathbb{R}} i \sum_{i=1}^m \underbrace{\max(0, 1 - y_i(g(x_i) + b))}_{l(g(x_i), y_i)} + \|g\|_{\mathcal{H}}^2$$

Classifier: $f(x) = \text{sign}(g(x) + b)$

- Using the representer theorem one can deduce that at optimality $g(\cdot) = \sum_{i=1}^m \gamma_i K(\cdot, x_i)$
- Plugging this into the optimization problem mentioned in the previous slide recovers the C-SVM formulation in the feature space.
- The Trace of the Kernel matrix plays an important role in the generalization abilities

Computation in feature spaces using dot products

$$K(x, z) = \langle \Phi(x), \Phi(z) \rangle \quad \|\Phi(x)\| = \sqrt{\langle \Phi(x), \Phi(x) \rangle} = \sqrt{K(x, x)}$$

Normalized Kernel

$$\hat{\Phi}(x) = \frac{\Phi(x)}{\|\Phi(x)\|} \quad \hat{K}(x, z) = \langle \hat{\Phi}(x), \hat{\Phi}(z) \rangle = \frac{K(x, z)}{\sqrt{K(x, x)K(z, z)}}$$

Distance between feature vectors

$$\begin{aligned} \|\Phi(x) - \Phi(z)\|^2 &= \langle \Phi(x) - \Phi(z), \Phi(x) - \Phi(z) \rangle \\ &= K(x, x) + K(z, z) - 2K(x, z) \end{aligned}$$

If Φ is normalized $K(x, x) = 1$ then $\|\Phi(x) - \Phi(z)\|^2 = 2 - 2K(x, z)$
 $K(x, z)$ can be understood as a measure of similarity between x and z .

Computation in feature spaces using dot products

Computing norms of linear combination of feature maps

$$\left\| \sum_{i=1}^m \alpha_i \Phi(x_i) \right\|^2 = \left\langle \sum_{i=1}^m \alpha_i \Phi(x_i), \sum_{j=1}^m \alpha_j \Phi(x_j) \right\rangle = \sum_{i=1, j=1}^m \alpha_i \alpha_j K(x_i, x_j)$$

Let $\bar{\Phi} = \frac{1}{m} \sum_{i=1}^m \Phi(x_i)$.

$$\|\bar{\Phi}\|^2 = \langle \bar{\Phi}, \bar{\Phi} \rangle = \frac{1}{m^2} \sum_{i,j=1}^m K(x_i, x_j)$$

$\Phi_{new}(x) = \Phi(x) - \bar{\Phi}$

$$K_{new}(x, z) = K(x, z) - \frac{1}{m} \sum_{i=1}^m K(x, x_i) - \frac{1}{m} \sum_{i=1}^m K(z, x_i) + \frac{1}{m^2} \sum_{i,j=1}^m K(x_i, x_j)$$

Computation in feature spaces using dot products

Assume that

$S_1 = \{x_1, \dots, x_{m_+}\}$ positive class

$S_2 = \{z_1, \dots, z_{m_-}\}$ negative class

Let $\bar{\Phi}_1$ be defined on S_1 and $\bar{\Phi}_2$ be defined on S_2

$$f(x) = \text{sign}(\|\Phi(x) - \bar{\Phi}_1\|^2 - \|\Phi(x) - \bar{\Phi}_2\|^2)$$

Computation in feature spaces using dot products

Assume that

$S_1 = \{x_1, \dots, x_{m_+}\}$ positive class

$S_2 = \{z_1, \dots, z_{m_-}\}$ negative class

Let $\bar{\Phi}_1$ be defined on S_1 and $\bar{\Phi}_2$ be defined on S_2

$$\begin{aligned} f(x) &= \text{sign}(\|\Phi(x) - \bar{\Phi}_1\|^2 - \|\Phi(x) - \bar{\Phi}_2\|^2) \\ &= \text{sign}\left(\frac{1}{m_+} \sum_{i=1}^{m_+} K(x, x_i) - \frac{1}{m_-} \sum_{i=1}^{m_-} K(x, z_i) - b\right) \end{aligned}$$

Summary

- reviewed the notion of dot products
- seen examples of how one can do computations in feature space
- reviewed kernel functions

PART 4: Kernels on Different Data types

Similarity between Strings

- Impetus from Biology
- *Insertion, Deletion and Substitution*

Let $x = \text{singapore}$ and $y = \text{ingpore}$

s i n g a p o r e
- i n g - p o r e

p o r e

Global alignment Needleman

Wunsch $O(|x||y|)$

is the largest common
subsequence.

Local alignment

Smith Waterman $O(|x||y|)$

Both algorithms gives a score, which is used heavily in computational biology.

Probabilistic approaches for modelling sequences

- Probabilistic approaches can be an interesting alternative (Durbin et al. 2000)
- HMMs have been used very successfully for modelling sequences
- A very special HMM model called **Profile HMM** was invented for probabilistically aligning **multiple sequences**

Spectrum Kernels(Leslie et al. 2002)

- Consider all k - length contiguous subsequences that it contains
- Create a feature map indexed by all possible substrings ("k-mers") from the alphabet of **amino** acids
- Dimension of feature space $|\Sigma|^k, \Sigma = 20$

```
s  i  n  g  a  p  o  r  e
s  i  n
   i  n  g
      n  g  a
         g  a  p
            a  p  o
               p  o  r
                  o  r  e
```

Spectrum kernel(Leslie et al. 2002)

Let $x = \text{singapore}$ and u be a substring of x of length 3(u is a 3-mer)

feature map

$$\phi_3(x) \in \mathbb{R}^l \quad \phi_3(x)_u = \# \text{ occurrences of } u \text{ in } x$$

$$\phi_3(x) \in \mathbb{R}^l \quad l = 26 \times 3$$

Spectrum kernel(Leslie et al. 2002)

Let $x = \text{singapore}$ and u be a substring of x of length 3(u is a 3-mer)

feature map

$$\phi_3(x) \in \mathbb{R}^l \quad \phi_3(x)_u = \# \text{ occurrences of } u \text{ in } x$$

$$\phi_3(x) \in \mathbb{R}^l \quad l = 26 \times 3$$

aaa	...	sin	ing	nga	gap	apo	por	ore	...	zzz
0	...	1	1	1	1	1	1	1	...	0

$$K(x, z) = \phi_3(x)^\top \phi_3(z)$$

Allowing for mismatches

For each u form a r mismatch neighbourhood, $N_r(u)$

$$N_{3,1}(\text{sin}) = \{\alpha \text{in}, s\beta \text{n}, \text{si}\gamma\}$$

where $\alpha, \beta, \gamma \in A$ and $\alpha \neq s, \beta \neq i, \gamma \neq n$

$$\phi_{3,1}(u) = (I_{v \in N_{3,1}(u)})_{v \in 3\text{-mer}}$$

$$\Phi_{3,1}^{mis} = \sum_{3\text{-mer } u \in x} \phi_{3,1}(u)$$

Computing $K(x, z)$ can be solved in $O\left(k^{r+1} |\Sigma|^r \underbrace{(|x| + |z|)}_{\text{linear}}\right)$

Variations on the theme (Leslie and Kuang, 2003)

- Introduce gaps
- probabilistic substitutions
- Wild cards

What is a Generative model

$$P(X = x|\theta)$$

Hidden Markov models

What is a Generative model

$$P(X = x|\theta)$$

Hidden Markov models

Given $D = \{x_1, \dots, x_m\}$ Maximum likelihood estimate is given by

$$\hat{\theta} = \operatorname{argmax}_{\theta} \sum_{i=1}^n \log P(x_i|\theta)$$

Fisher kernel (Jaakkola and Haussler 1999)

Fisher information matrix

$$F(\theta) = \int_x U(x, \theta) U(x, \theta)^\top P(x|\theta) dx$$

$$U(x, \theta) = \nabla_\theta \log P(x|\theta)$$

Fisher Kernel

$$K_f(x, z) = U(x, \hat{\theta})^\top F^{-1}(\hat{\theta}) U(z, \hat{\theta})$$

Fisher kernel (Jaakkola and Haussler 1999)

Computing F is hard but can be approximated $F \approx I$
If $x \sim$ Exponential family

$$\log P(x|\theta) = s(x) + \Psi(\theta)$$

$s(x)$ sufficient statistics

$$K_f(x, z) = s(x)^\top s(z) + c$$

Graphs play a very important role in Computational Biology

- Ideal to model molecules
- Protein-protein interaction networks
- metabolic networks

See (Vishwanathan et al. 2010, Gartner et al. 2003, Kashima et al. 2003)

Diffusion Kernels(Kondor and Lafferty 2002)

Let $\mathcal{X} = \{1, \dots, m\}$ and there be some associated edges between them. The adjacency matrix of the resulting graph be A .

Diffusion Kernel

$$K = \lim_{s \rightarrow \infty} \left(I + \frac{\beta}{s} H \right)^s$$

$H = A - D$, where D is diagonal with $d_{ii} = \sum_j a_{ij}$. K is positive definite and Symmetric. Computation is $O(m^3)$

Diffusion Kernels(Kondor and Lafferty 2002)

Let $\mathcal{X} = \{1, \dots, m\}$ and there be some associated edges between them. The adjacency matrix of the resulting graph be A .

Diffusion Kernel

$$K = \lim_{s \rightarrow \infty} \left(I + \frac{\beta}{s} H \right)^s$$

$H = A - D$, where D is diagonal with $d_{ii} = \sum_j a_{ij}$. K is positive definite and Symmetric. Computation is $O(m^3)$

$$\lim_{s \rightarrow \infty} \left(1 + \frac{\beta}{s} x \right)^s = e^{\beta x}$$

$K = e^{\beta H} = \sum_{i=1}^m v_i e^{\beta \lambda_i} v_i$ where (λ_i, v_i) are the (eigen-value, eigen-vector) of H

Diffusion Kernels(Kondor and Lafferty 2002)

Sometimes can be computed in closed form for special graphs e.g.
complete graphs

$$K(i,j) = \begin{cases} \frac{1+(m-1)e^{-m\beta}}{m} & i = j \\ \frac{1-e^{-m\beta}}{m} & i \neq j \end{cases}$$

Has a very interesting analogue with Diffusion equation in physics.

From Similarity measures to Kernels

In many applications Similarity measures are available. Let $s : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a similarity measure.

Kernel from Empirical Kernel Map

$$K(x, z) = \phi(x)^\top \phi(z) = \sum_{i=1}^l s(x, t_i) s(z, t_i)$$

$\phi(x, t_i) = [s(x, t_1) \dots, s(x, t_l)]^\top$ (Tsuda, 1999) is called the empirical kernel map. The elements t_1, \dots, t_l are pre-specified and are called templates

Found useful in remote protein Homology detection (Liao and Noble, 2002).

Other domains

- Kernels on Probability distributions (See Kondor et al. 2004)
- An interesting set of Kernels based on Binet-Cauchy theorem was proposed by (Vishwanathan and Smola, 2004) useful for dynamical systems
- Exploiting Semigroups, Kernels on measures was derived in (Cuturi et al. 2005)

PART 5: Advanced Topics

Recap of SVMs

- On a dataset $D = \{(x_i, y_i) | i = 1, \dots, m\}$ SVMs solve the following problem

$$\Gamma(K) = \max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \alpha^\top YKY\alpha \quad (1)$$

$$0 \leq \alpha_i \leq C \quad \sum_i \alpha_i y_i = 0 \quad (2)$$

where $K_{ij} = k(x_i, x_j)$ is the kernel function defined on examples x_i and x_j .

- The final classifier is $y = \text{sign}(\sum_i \alpha_i y_i k(x, x_i) + b)$

Recap of SVMs

- Does not scale well
- The function $\Gamma(K)$ is a pointwise maximum of a set of functions and hence **is convex**
- If the maximization in α is not unique then $\Gamma(K)$ is not differentiable.
- $\Gamma(K)$ maynot be differentiable, but *subgradients* exist.
- Let us relax the problem a little and say that $\mu_i \geq 0$

On the Problem of Kernel Design

- Learn an optimal Kernel from a library of kernel functions
- Learn an optimal kernel from several **similarity** functions
- Find a *good* classifier when K is **uncertain**
- All three problems gives rise to extremely interesting optimization problems

MKL is a Semi-definite Programming problem

$$\begin{aligned} \min_z \quad & c^\top z \\ \text{s.t.} \quad & F(z) = \sum_{i=1}^l z_i F_i \succeq 0 \\ & Bz = d \end{aligned}$$

$z \in \mathbb{R}^l$ and $F_i = F_i^\top \in \mathbb{R}^{m \times m}$.

$F(z)$ is positive semidefinite

Instance of a convex optimization problem. Can be solved by interior point methods

Learning a linear combination of Multiple Kernels

Let $\{K_1, \dots, K_n\}$ be a given library of kernels. Given a training set of m examples, each $K_i = K_i^\top \in \mathbb{R}^{m \times m}$

Find a good classifier using a kernel function of the form $\sum_{i=1}^l \mu_i K_i$ such that $\Gamma(K)$ is minimized (Lanckriet et al. 2004) The problem

requires $K \succeq 0$ and $\text{trace}(K) = c$ (for theoretical guarantees)

MKL formulation

SDP formulation

$$\min_{\mu, t, \lambda, \mathbf{v} > 0} \quad t \quad (3)$$

$$\begin{pmatrix} \sum_{i=1}^l \mu_i Y K_i Y^\top & e + \mathbf{v} + \lambda y \\ e + \mathbf{v} + \lambda y & t \end{pmatrix} \succeq 0 \quad (4)$$

$$\sum_{i=1}^n \mu_i K_i \succeq 0 \quad (5)$$

Reformulation of MKL

The SDP problem can be recast as QCQPs

$$\max_{\alpha, t} \quad \alpha^\top e - ct \quad (6)$$

$$s.t. \quad \alpha^\top YK_i Y \alpha \leq r_i t \quad i = 1, \dots, l \quad (7)$$

$$\alpha^\top y = 0 \quad 0 \leq \alpha \leq C \quad (8)$$

where $r_i = \text{trace}(K_i)$

QCQPs are instances of SOCPs

$$\min_z \quad c^\top z \quad (9)$$

$$\|A_i z + b_i\|_2 \leq c_i^\top z + d_i \quad (10)$$

where $A_i \in \mathbb{R}^{n_i \times l}$, $b_i, c_i, c, z \in \mathbb{R}^l$, $d_i \in \mathbb{R}$

Equivalence with Block L1 regularization

Bach et al. (2004) showed that the QCQP formulation is equivalent to

$$\min_{w,b,\xi} \quad \frac{1}{2} \left(\underbrace{\sum_{i=1}^l d_i \|w_i\|}_{\text{Block L1}} \right)^2 + C \sum_{i=1}^m \xi_i \quad (11)$$

$$s.t. \quad y_i (\sum_j w_j^\top \phi_j(x_i) + b) \geq 1 - \xi_i \quad \forall i = \{1, \dots, m\} \quad \xi_i \geq 0 \quad (12)$$

for proper choice of d_i

Block L1 norm promotes sparsity i.e. most of $\mu_i = 0$

- Scalable solutions to MKL
- Non-sparse formulations

Efficient algorithms for MKL

A trick

Let $\gamma \in \mathbb{B}_n = \{\gamma \in \mathbb{R}^n \mid \gamma_i \geq 0, \sum_{i=1}^n \gamma_i = 1\}$ For any $a_i \in \mathbb{R}, i = 1, \dots, n$

$$\left(\sum_{i=1}^n |a_i|\right)^2 \leq \sum_{i=1}^n \frac{a_i^2}{\gamma_i}$$

This implies that

$$\left(\sum_{i=1}^n \|w_i\|\right)^2 \leq \sum_{i=1}^n \frac{1}{\gamma_i} \|w_i\|^2$$

where γ lies in a probability simplex

Can be helpful in reformulating the L1 formulation.

Solving MKL by reusing SVM solvers

The following problem is equivalent to the Block L1 formulation (Rakotomamonjy et al. 2007)

$$S_m = \{\alpha \mid 0 \leq \alpha_i \leq C, \alpha^\top y = 0\} \text{ and } \mathbb{B} = \{\mu \mid 0 \leq \mu_i, \sum_{i=1}^l \mu_i = 1\}$$

$$\min_{\mu \in \mathbb{B}} J(\mu) \quad \left(= \max_{\alpha \in S_m} \alpha^\top e - \frac{1}{2} \sum_{i=1}^l \mu_i \alpha^\top Y K_i Y \alpha \right) \quad (13)$$

A gradient descent algorithm-iteration

- 1.) Solve SVM problem with Kernel $K = \sum_{i=1}^l \mu_i K_i$
- 2.) Differentiate J w.r.t μ and update μ

See also Sonnenberg et al. 2006

Motivation for Non-sparse MKL setting

- In many situations sparsity may be suboptimal
- (Kloft et al. 2009) one can pursue general norms
- (Nilsback and Zisserman, 2006) points out that in object categorization tasks, employing a few of the **feature descriptors** or employing a canonical combination of them often leads to sub-optimal solutions.
- Simply using a sparsity criterion or a non-sparsity criterion is not going to work

Motivation for Non-sparse MKL setting

- In object categorization (Nilsback and Zisserman, 2006) study several descriptors related to **shape**, **color** and **texture**. They studied a total of seven descriptors to classify flowers.
- MKL learning is of extreme importance to object categorization tasks.
- Goal is to design a suitable **regularization term** which could select all the descriptors and yet be **able to generalize well**

Hierarchical feature space

Training dataset: $\mathcal{D} = \{(\mathbf{x}_i, y_i), i = 1, \dots, m \mid \mathbf{x}_i \in \mathcal{X}, y_i \in \{-1, 1\}\}$.
Let each individual example be described by a concatenation of features

$$\mathbf{x} = \left\{ \begin{array}{l} \left\{ \begin{array}{l} \mathbf{x}_{j1} \\ \vdots \\ \mathbf{x}_{jk} \\ \vdots \\ \mathbf{x}_{jn_j} \end{array} \right. \\ \vdots \end{array} \right.$$

where $\mathbf{x}_{jk} \in \mathbb{R}^{d_{jk}}$ is a vector of attributes which correspond to the j, k th block. Let $K_{jl}(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_{1jl}^\top \mathbf{x}_{2jl}$

Regularizing \mathbf{w}

$$y = \text{sign} \left(\sum_{j=1}^n \sum_{k=1}^{n_j} w_{jk}^\top \mathbf{x}_{jk} + b \right)$$

- Learning objective: $J(\mathbf{w}) = \Omega(\mathbf{w}) + cL(\mathcal{D})$

A new regularization term

$$\Omega(\mathbf{w}) = \frac{1}{2} \left\{ \sum_{j=1}^n \left\{ \sum_{k=1}^{n_j} \|\mathbf{w}_{jk}\|_2 \right\}^{2q} \right\}^{\frac{1}{q}}$$

for $q \geq 1$

- Block L_1 norm is given by $q = 1$, and $n = 1$.

A general formulation

We consider the following formulation

$$\begin{aligned} \min_{\mathbf{w}_{jk}, b, \xi_i} \quad & \frac{1}{2} \left[\sum_j \left(\sum_{k=1}^{n_j} \|\mathbf{w}_{jk}\|_2 \right)^{2q} \right]^{\frac{1}{q}} + C \sum_i \xi_i \\ \text{s.t.} \quad & y_i \left(\sum_{j=1}^n \sum_{k=1}^{n_j} \mathbf{w}_{jk}^\top \mathbf{x}_{jki} - b \right) \geq 1 - \xi_i, \quad \xi_i \geq 0 \quad \forall i \end{aligned} \quad (14)$$

If each group contains just one element then one generalizes the MKL setting

We present a convergent algorithm, $O(m^2 n^2 \ln n_t / \epsilon^2)$,
 m = no of examples, n = no of groups, $n_t = \max_j n_j$ which solves such a problem for arbitrary $q \geq 1$.

Learning kernels from similarity functions

- Often similarity functions are more readily available than kernels.
- **Sequence Similarity** and **Structure Similarity** is readily available for **Protein structures**
- Classifying protein structures is an important problem where both sequence and structure are important

Problem Definition

Given: a set of m training examples with similarity matrices $S_1, S_2, \dots, S_n \in \mathbb{R}^{m \times m}$, and class labels $\mathbf{y} \in \{+1, -1\}^m$.

Task: learn a kernel matrix K , such that SVM performs well on test data.

$$\min_{K \succeq 0} \Gamma(K) + \rho \|K - S\|_F^2. \quad (15)$$

- It is a minimax problem

$$\min_K \max_{\alpha} f(\alpha, K)$$

- First studied in (Luss & Aspremont, 2007) and solved using **Analytic center cutting plane** method.
- (Chen et al et al, 2008) formulated the above problem as **Quadratically Constrained Linear Program**.
- (Chen et al. 2009) used an alternate loss function $\|K - S\|_F$ which led to **Second Order Cone Program**.

Related work

- Exploits strong duality

$$\min_K \max_{\alpha} f(K, \alpha) = \max_{\alpha} \min_K f(K, \alpha)$$

- For Frobenius Norm Loss one could solve $\min_K f(K, \alpha)$ in a closed form

$$K^*(\alpha) = \left(S + \frac{1}{4\rho} (Y\alpha\alpha^T Y) \right)_+$$

(Luss & Aspremont 2007)

- One can solve the resulting problem

$$\max_{\alpha} f(K^*(\alpha), \alpha)$$

by ACCP procedure

- Will not apply to other losses

Single Kernel based Formulation

$$\min_K \left[\Gamma(K) + \rho \sum_{l=1}^n L_l(K - S_l) \right],$$

$$s.t. \quad K \succeq 0, \quad \text{trace}(K) = \tau.$$

$L_l(\cdot)$ is a convex, sub-differentiable loss function.

Examples:

$$1] L_l(K - S_l) = \sum_i \sum_j |K(i,j) - S_l(i,j)|$$

$$2] L_l(K - S_l) = \|K - S_l\|_F$$

$$3] L_l(K - S_l) = \sum_i \sum_j [K(i,j) - S_l(i,j)]^2$$

Recall: Multiple Kernel Learning

$$\Gamma(K_1, \dots, K_n) = \max_{\gamma \in \Delta} \max_{\alpha \in A_m} \left[\alpha^\top \mathbf{1} - \sum_{l=1}^n \frac{\alpha^\top Y K_l Y \alpha}{2 \gamma_l} \right],$$

$$\mathbb{B}_n = \{ \gamma \in \mathbb{R}^n \mid \gamma \geq 0, \gamma^\top \mathbf{1} = 1 \}.$$

Optimal kernel:
$$K^* = \sum_{l=1}^n \frac{1}{\gamma_l^*} K_l$$

Multiple Kernel based Formulation

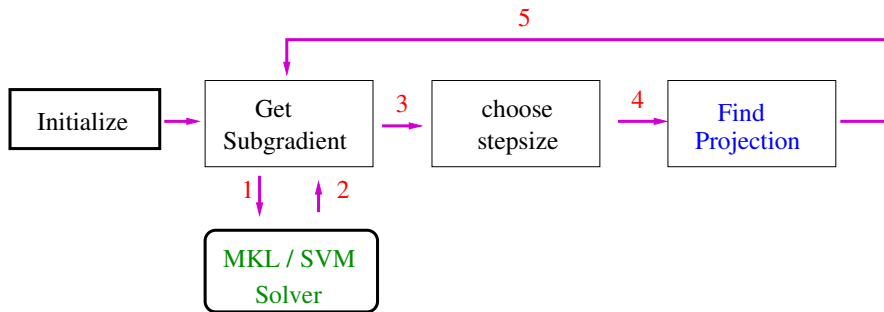
$$\min_{K_1, \dots, K_n} \left[\Gamma(K_1, \dots, K_n) + \rho \sum_{l=1}^n L_l(K_l - S_l) \right],$$

$$s.t. \quad K_l \succeq 0, \text{ trace}(K_l) = c, \quad l = 1, \dots, n.$$

Ideas behind the algorithm

- Solve the inner maximization over α for a fixed Kernel or Kernels by solving an SVM or MKL algorithm
- Solve the outer minimization by SubGradient projection.
- We choose Mirror Descent for the Subgradient projection step

Mirror Descent based Algorithm



MD Step 1: Subgradient Computation

- Obtain (α^*, γ^*) by solving the no n-sparse MKL.
- Calculate subgradient of Γ as

$$\Gamma'(K_1, \dots, K_n) = (\Gamma'_1, \dots, \Gamma'_n),$$

$$\Gamma'_l := -\frac{1}{2\gamma_l^*} Y \alpha^* \alpha^{*\top} Y.$$

- Compute L'_l - a subgradient of L_l .

MD Step 2: Compute Projection

- Eigen decomposition:

$$\begin{aligned} \Gamma'_l(K_1^{(t)}, \dots, K_n^{(t)}) + \rho L'_l(K_l^{(t)} - S_l) \\ = V_l^{(t)} \text{diag}([d_{1,l}^{(t)} \dots d_{n,l}^{(t)}]) V_l^{(t)\top}. \end{aligned}$$

- $\lambda_{i,l}^{(t+1)} := \frac{c \lambda_{i,l}^{(t)} \exp(-\eta_t d_{i,l}^{(t)})}{\sum_{j=1}^n \lambda_{j,l}^{(t)} \exp(-\eta_t d_{j,m}^{(t)})}, i = 1, \dots, n.$
- $K_l^{(t+1)} := V_l^{(t)} \text{diag}([\lambda_{1,l}^{(t+1)} \dots \lambda_{n,l}^{(t+1)}]) V_l^{(t)\top}.$

Algorithmic Convergence

$F^{(t)}$: objective function value at t -th iteration.

F^* : optimal objective value.

$Lip(F)$: Lipschitz constant of objective function .

If the algorithm is initialized with $K_l^{(1)} = \frac{c}{m}I, \forall n$ and the stepsizes are chosen as $\eta_t = \frac{1}{Lip(F)} \sqrt{\frac{2 \log m}{nt}}$

then $\min_{1 \leq t \leq T} F^{(t)} - F^ \leq cnLip(F) \sqrt{\frac{2 \log m}{T}}$.*

Algorithmic Complexity

$O\left(\frac{n^2 \log m}{\epsilon^2}\right)$ iterations to reach ϵ -accurate solution.

At every iteration key computations are:

- Eigen decomposition of n matrices of dimension $m \times m$.
- Solving SVM / MKL with m training examples and n kernels.

Contributions

- Note that $\max_{\alpha} f(K, \alpha)$ is equivalent to solving an SVM. Unlike the state of the art we intend to solve the original problem

$$\min_K \max_{\alpha} f(K, \alpha)$$

- proposed three formulations which could be solved in an iterative fashion
- Each iteration is no expensive than solving an MKL/SVM
- Leads to scalable solution
- State of the art seems to be specific for the frobenius loss and cannot handle the general case
- Previous work could not handle multiple similarity matrices,

References

Kernel methods in Computational Biology

Scholkopf et al. 2004

Kernel methods for Pattern Analysis

John Shawe Taylor and N. Cristianini

Learning with Kernels

Scholkopf and Smola 2002