

Multiple Kernels and Reproducing Kernel Hilbert Spaces

Lecturer: Michael I. Jordan

Scribes: Blaine Nelson and Sierra Boyd

1 Convex Programming

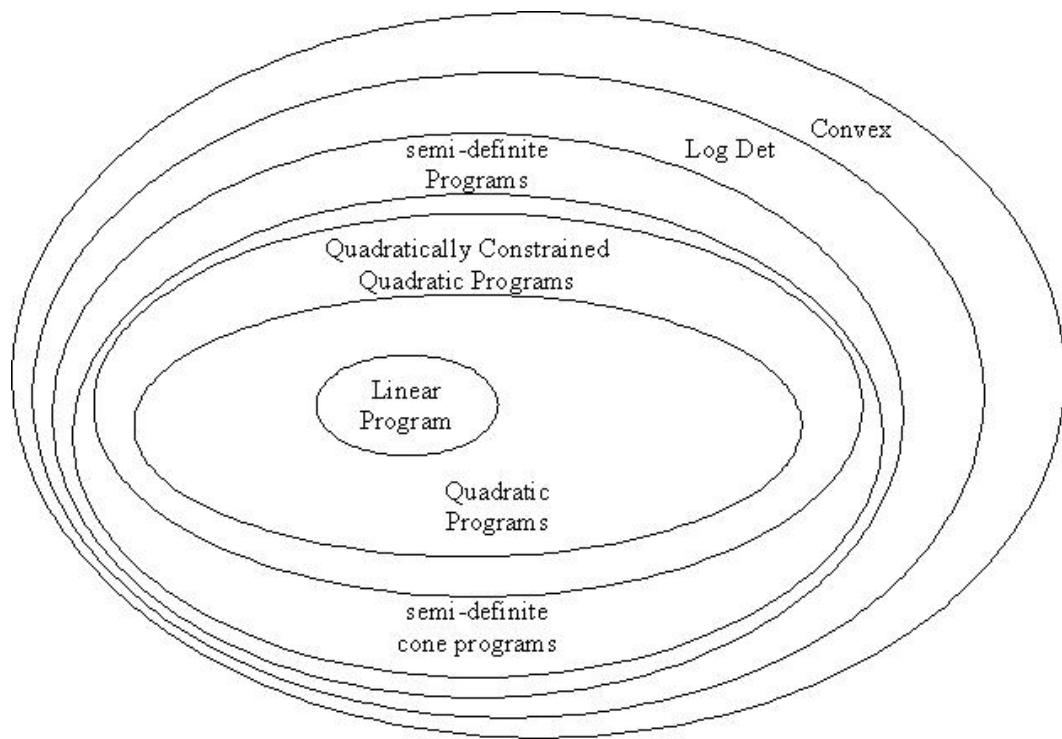


Figure 1: Diagram of Convex Programming

2 Multiple Kernels

We continue our discussion on multiple kernels from the previous lecture. Recall,

$$\begin{aligned} K &= \sum_i \mu_i K_i \\ K &\succeq 0 \\ \text{tr}(K) &\leq c \end{aligned}$$

Given data pairs $\{(x_i, y_i)\}_{i=1}^n$ ($x_i \in \mathcal{X}$ and $y_i \in \{-1, 1\}$), let $G(\mathbf{K}) = \text{diag}(Y) \cdot K \cdot \text{diag}(Y)$ where

$$\text{diag}(Y) = \begin{pmatrix} y_1 & & & \\ & y_2 & \mathbf{0} & \\ & & \ddots & \\ & & & y_n \end{pmatrix}$$

Thus, the (i,j)-th element of $G(\mathbf{K})$ is given by $G(\mathbf{K})_{ij} = y_i y_j K(x_i, x_j)$

The hard-margin SVM associated with $G(\mathbf{K})$ is represented by the following program:

$$\begin{aligned} \max \quad & 2\alpha^T e - \alpha^T G(\mathbf{K}) \alpha \\ \text{s.t.} \quad & \alpha^T y = 0 \quad \alpha_i \geq 0 \end{aligned}$$

$$e = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

Moreover, the soft-margin SVM can be represented by a similar program where the last constraint is boxed $0 \leq \alpha_i \leq C$.

How can we optimize over α and μ ?

Theorem 1

Multiple Kernel Optimization Theorem: *The problem of optimizing*

$$\begin{aligned} \max_{\alpha} \quad & 2\alpha^T e - \alpha^T (G(\mathbf{K}) + \tau \mathbf{I}) \alpha \\ \text{s.t.} \quad & \alpha^T y \geq 0 \quad 0 \leq \alpha_i \leq C \end{aligned}$$

with respect to $\{\mu_i\}$ where τ is a regularization parameter reduces to:

$$\begin{aligned} \min_{\{\mu_i\}, t, \lambda, \nu, \delta} \quad & t \\ \text{s.t.} \quad & \text{tr}(\mathbf{K}) = 0 \quad \mathbf{K} = \sum_i \mu_i K_i \\ & \begin{pmatrix} G(\mathbf{K}) + \lambda \mathbf{I} & e + \nu - \delta + \lambda y \\ (e + \nu - \delta + \lambda y)^T & t - 2C\delta^T e \end{pmatrix} \succeq 0 \\ & \nu \geq 0 \quad \delta \geq 0 \end{aligned}$$

which is a Semi-Definite Program (SDP).

Proof. Optimizing with respect to $\{\mu_i\}$ is equivalent to the following program:

$$\begin{aligned} \min_{\{\mu_i\}} \max_{\alpha} \quad & 2\alpha^T e - \alpha^T (G(\mathbf{K}) + \tau \mathbf{I}) \alpha \\ \equiv \min_{\{\mu_i\}, t} \quad & t \geq \max_{\alpha} (2\alpha^T e - \alpha^T (G(\mathbf{K}) + \tau \mathbf{I}) \alpha) \end{aligned}$$

Taking the dual of this gives the Lagrangian,

$$L(\alpha, \nu, \lambda, \delta) = 2\alpha^T e - \alpha^T (G(\mathbf{K}) + \tau \mathbf{I}) \alpha + 2\nu^T \alpha + 2\lambda y^T \alpha + 2\delta^T (Ce - \alpha)$$

Now, we maximize over α to find the dual problem. This is done by setting $\frac{\partial L}{\partial \alpha}$ to 0. Thus we find that $\alpha = (G(\mathbf{K}) + \tau \mathbf{I})^{-1} (e + \nu - \delta + \lambda y)$. This yields the following dual problem:

$$\min_{\nu \geq 0, \delta \geq 0, \lambda} \underbrace{(e + \nu - \delta + \lambda y)^T (G(\mathbf{K}) + \tau \mathbf{I})^{-1} (e + \nu - \delta + \lambda y) + 2C\delta^T e}_{\beta}$$

Hence the primal statement is less than or equal to t if and only if $\exists \nu, \delta, \lambda$ such that $\beta \leq t$.

But, from the Schur Complement Theorem, we know that,

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix} \succeq 0 \Leftrightarrow \begin{matrix} A \succeq 0 \\ D - CA^{-1}B \succeq 0 \end{matrix}$$

But $D - CA^{-1}B$ has the exact same form as $-\beta$. Moreover, from the problem statement, $(G(\mathbf{K}) + \tau \mathbf{I}) \succeq 0$. Thus, $(e + \nu - \delta + \lambda y)^T (G(\mathbf{K}) + \tau \mathbf{I})^{-1} (e + \nu - \delta + \lambda y) + 2C\delta^T e \leq 0$ if and only if,

$$\begin{pmatrix} G(\mathbf{K}) + \tau I & e + \nu - \delta + \lambda y \\ (e + \nu - \delta + \lambda y)^T & t - 2C\delta^T e \end{pmatrix} \succeq 0$$

Thus, this condition implies that $\beta \leq t$, which is equivalent to the original min-max problem. □

3 Hilbert Space

A Hilbert space is essentially an Euclidean space, but a Euclidean space that may be infinite-dimensional. It is a vector space (i.e., is closed under addition and scalar multiplication, obeys the distributive and associative laws, etc.). It is also endowed with an inner product $\langle \cdot, \cdot \rangle$; a bilinear form obeying the following conditions:

$$\begin{aligned} \langle x + y, z \rangle &= \langle x, z \rangle + \langle y, z \rangle \\ \langle \alpha x, y \rangle &= \alpha \langle x, y \rangle \\ \langle x, y \rangle &= \langle y, x \rangle \\ \langle x, x \rangle &\geq 0 \\ \langle x, x \rangle = 0 &\implies x = 0 \end{aligned}$$

From $\langle \cdot, \cdot \rangle$ we get a norm $\| \cdot \|$ via $\| x \| \triangleq \langle x, x \rangle^{1/2}$. This norm allows us to define notions of convergence. Adding all limit points of Cauchy sequences (the distance between elements x_i and x_{i+1} of the sequence becomes small as $i \rightarrow \infty$) to our space yields a Hilbert space—a *complete* inner product space.

Theorem 2

Cauchy-Schwartz:

$$\langle x, y \rangle \leq \| x \| \| y \|$$

is easily proved for any Hilbert space.

3.1 Examples

- R^n : $\langle x, y \rangle = x^T y$
- L_2 : $\langle x, y \rangle = \int x(t)y(t)dt$
- l_2 : $\langle x, y \rangle = \sum_{i=1}^{\infty} x_i y_i$

4 Reproducing Kernel Hilbert Spaces

The Hilbert space L_2 is too “big” for our purposes, containing too many non-smooth functions. One approach to obtaining restricted, smooth spaces is the Reproducing Kernel Hilbert Space (RKHS) approach. A RKHS is “smaller” than a general Hilbert space.

Given a kernel $k(x, x')$, we will construct a Hilbert space such that k defines an inner product in that space. First define the *Gram matrix*. Given points $x_1, x_2, x_3, \dots, x_n$, define:

$$K_{ij} = k(x_i, x_j)$$

We say that the kernel k is *positive definite* (p.d.) if its Gram matrix is positive definite for all x_1, x_2, \dots, x_n . Similarly, we define positive semidefinite (p.s.d.), negative definite, and negative semidefinite by the properties of the Gram matrix.

The Cauchy-Schwartz inequality holds for p.s.d. kernels:

$$k(x_1, x_2)^2 \leq k(x_1, x_1)k(x_2, x_2)$$

Proof: Form a Gram matrix of the two points x_1 and x_2 :

$$K = \begin{pmatrix} k(x_1, x_1) & k(x_1, x_2) \\ k(x_2, x_1) & k(x_2, x_2) \end{pmatrix}$$

For K to be positive semidefinite as a matrix, the determinant of K must be nonnegative:

$$\implies k(x_1, x_1)k(x_2, x_2) - k(x_2, x_1)k(x_1, x_2) \geq 0,$$

which implies Cauchy-Schwartz.

Define the following *reproducing kernel map* for a kernel $k(\cdot, \cdot)$:

$$\Phi : x \longrightarrow k(\cdot, x).$$

I.e., to each point x in the original space we associate a function $k(\cdot, x)$.

Example: Gaussian kernel. Each point x maps to a Gaussian centered at that point as depicted in Fig. 2. Intuitively this captures the similarity of x to all other points.

We now construct a vector space containing all linear combinations of the functions $k(\cdot, x)$:

$$f(\cdot) = \sum_{i=1}^m \alpha_i k(\cdot, x_i).$$

for arbitrary $\{x_i \in \mathcal{X}\}_{i=1}^m$. This will be used to construct our RKHS H .

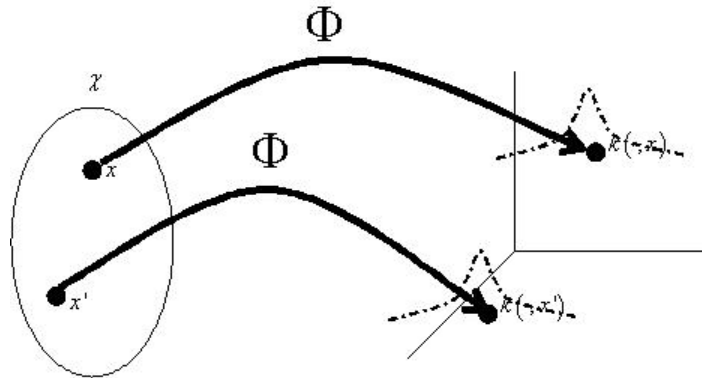


Figure 2: The mapping of an input space to a gaussian feature space.

We now define an inner product. Let $f(\cdot) = \sum_i \alpha_i k(\cdot, x_i)$ and $g(\cdot) = \sum_j \beta_j k(\cdot, x'_j)$ such that $f, g \in H$, and define:

$$\langle f, g \rangle \triangleq \sum_{i=1}^m \sum_{j=1}^{m'} \alpha_i \beta_j k(x_i, x'_j)$$

We need to verify that this in fact defines an inner product. Symmetry is obvious: $\langle f, g \rangle = \langle g, f \rangle$. Linearity is easy to show. In the next lecture we will establish the key property: $\langle f, f \rangle = 0 \implies f = 0$.

Kernels are analogs of Dirac delta functions. Consider L_2 (which is not a RKHS). We have:

$$f(x) = \int f(t) \delta(t, x) dt,$$

where $\delta(t, x)$ is the Dirac delta function. The Dirac delta function is the representer of evaluation for L_2 , but of course it is not itself in L_2 . (Which is consistent with the fact that L_2 is not a RKHS).